

L_p -norm Sauer-Shelah Lemma for Margin Multi-category Classifiers

Yann Guermeur

LORIA-CNRS

Campus Scientifique, BP 239

54506 Vandœuvre-lès-Nancy Cedex, France

(e-mail: Yann.Guermeur@loria.fr)

September 25, 2016

Running Title: Sauer-Shelah Lemma for Margin Multi-category Classifiers

Keywords: margin multi-category classifiers, guaranteed risks, ϵ -entropy, γ -dimension, generalized Sauer-Shelah lemmas

Mathematics Subject Classification: 68Q32, 62H30

Abstract

In the framework of agnostic learning, one of the main open problems of the theory of multi-category pattern classification is the characterization of the way the complexity varies with the number C of categories. More precisely, if the classifier is characterized only through minimal learnability hypotheses, then the optimal dependency on C that an upper bound on the probability of error should exhibit is unknown. We consider margin classifiers. They are based on classes of vector-valued functions with one component function per category, and the classes of component functions are uniform Glivenko-Cantelli classes. For these classifiers, an L_p -norm Sauer-Shelah lemma is established. It is then used to derive guaranteed risks in the L_∞ and L_2 -norms. These bounds improve over the state-of-the-art ones with respect to their dependency on C , which is sublinear.

1 Introduction

During a long period, the theory of multi-category pattern classification was considered as a topic of limited importance. Two connected reasons can be put forward to explain this phenomenon. On the one hand, the theory dedicated to dichotomies was making rapid strides, on the other hand, decomposition methods were seen as efficient solutions to tackle polytomies. An obvious drawback of this line of reasoning is to neglect the specificities of the multi-category case, such as the dependency of the complexity of the task on the number C of categories. In recent years, several studies addressed this question, by deriving upper bounds on the probability of error of multi-category classifiers, especially margin ones. However, most of these *guaranteed risks* were dedicated to specific families of classifiers, let them be kernel machines [36, 24], neural networks [2], decision trees [23] or nearest neighbors classifiers [22]. This article deals with margin classifiers. They are based on classes of vector-valued functions with one component function per category, and the classes of component functions are uniform Glivenko-Cantelli classes. For these classifiers, an L_p -norm Sauer-Shelah lemma is established. It is then used to derive guaranteed risks in the L_∞ and L_2 -norms. These bounds improve over the state-of-the-art ones with respect to their dependency on C , which is sublinear. Thus, they pave the way for the characterization of the optimal dependency on C that could be obtained in the framework of agnostic learning, under minimal learnability/measurability hypotheses regarding the classes of functions involved.

The organization of the paper is as follows. Section 2 deals with the theoretical frame-

work and the margin multi-category classifiers. Section 3 is devoted to the derivation of the L_p -norm Sauer-Shelah lemma. The bound based on the L_∞ -norm and that based on the L_2 -norm are respectively established in Section 4 and Section 5. At last, we draw conclusions and outline our ongoing research in Section 6. To make reading easier, basic results from the literature and technical lemmas have been gathered in appendix.

2 Margin multi-category classifiers

The theoretical framework for the margin multi-category classifiers has been introduced in [16]. It is summarized below.

2.1 Theoretical framework

We consider the case of C -category pattern classification problems [12] with $C \in \mathbb{N} \setminus \llbracket 0, 2 \rrbracket$. Each object is represented by its description $x \in \mathcal{X}$ and the set \mathcal{Y} of the categories y can be identified with the set of indices of the categories: $\llbracket 1, C \rrbracket$. We assume that $(\mathcal{X}, \mathcal{A}_\mathcal{X})$ and $(\mathcal{Y}, \mathcal{A}_\mathcal{Y})$ are measurable spaces and denote by $\mathcal{A}_\mathcal{X} \otimes \mathcal{A}_\mathcal{Y}$ the tensor-product sigma algebra on the Cartesian product $\mathcal{X} \times \mathcal{Y}$. We make the hypothesis that the link between descriptions and categories can be characterized by an unknown probability measure P on the measurable space $(\mathcal{X} \times \mathcal{Y}, \mathcal{A}_\mathcal{X} \otimes \mathcal{A}_\mathcal{Y})$. Let $Z = (X, Y)$ be a random pair with values in $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, distributed according to P . The single knowledge source on P available is an m -sample $\mathbf{Z}_m = (Z_i)_{1 \leq i \leq m} = ((X_i, Y_i))_{1 \leq i \leq m}$ made up of independent copies of Z (in short $\mathbf{Z}_m \sim P^m$). The theoretical framework is thus that of *agnostic learning* [19]. To simplify reasoning, in the sequel, the hypothesis $m > C$ is made.

We add an hypothesis to that framework: the fact that the classifiers considered are based on classes of vector-valued functions with one component function per category, and the classes of component functions are *uniform Glivenko-Cantelli*. The definition of this property calls for the introduction of an intermediate definition.

Definition 1 (Empirical probability measure) *Let $(\mathcal{T}, \mathcal{A}_\mathcal{T})$ be a measurable space and let T be a random variable with values in \mathcal{T} , distributed according to a probability measure P_T on $(\mathcal{T}, \mathcal{A}_\mathcal{T})$. For $n \in \mathbb{N}^*$, let $\mathbf{T}_n = (T_i)_{1 \leq i \leq n}$ be an n -sample made up of independent copies of T . The empirical measure supported on this sample, $P_{\mathbf{T}_n}$, is given by*

$$P_{\mathbf{T}_n} = \frac{1}{n} \sum_{i=1}^n \delta_{T_i},$$

where δ_{T_i} denotes the Dirac measure centered on T_i .

Definition 2 (Uniform Glivenko-Cantelli class [15]) *Let the probability measures P_T and P_{T_n} be defined as in Definition 1. Let \mathcal{F} be a class of measurable functions on \mathcal{T} . Then \mathcal{F} is a uniform Glivenko-Cantelli class if for every $\epsilon \in \mathbb{R}_+^*$,*

$$\lim_{n \rightarrow +\infty} \sup_{P_T} \mathbb{P} \left(\sup_{n' \geq n} \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{T' \sim P_{T_{n'}}} [f(T')] - \mathbb{E}_{T \sim P_T} [f(T)] \right| > \epsilon \right) = 0,$$

where \mathbb{P} denotes the infinite product measure P_T^∞ .

Henceforth, we shall refer to uniform Glivenko-Cantelli classes by the abbreviation *GC classes*. GC classes must be uniformly bounded up to additive constants (see for instance Proposition 4 in [15]). For notational convenience, we replace this property by a stronger one: the vector-valued functions take their values in a hypercube of \mathbb{R}^C . The definition of a margin multi-category classifier is thus the following one.

Definition 3 (Margin multi-category classifiers) *Let $\mathcal{G} = \prod_{k=1}^C \mathcal{G}_k$ be a class of functions from \mathcal{X} into $[-M_{\mathcal{G}}, M_{\mathcal{G}}]^C$ with $M_{\mathcal{G}} \in [1, +\infty)$. The classes \mathcal{G}_k of component functions are supposed to be GC classes. For each function $g = (g_k)_{1 \leq k \leq C} \in \mathcal{G}$, a margin multi-category classifier on \mathcal{X} is obtained by application of the operator dr from \mathcal{G} into $(\mathcal{Y} \cup \{*\})^{\mathcal{X}}$ named decision rule and defined as follows:*

$$\forall x \in \mathcal{X}, \quad \begin{cases} |\arg\max_{1 \leq k \leq C} g_k(x)| = 1 \implies dr_g(x) = \arg\max_{1 \leq k \leq C} g_k(x) \\ |\arg\max_{1 \leq k \leq C} g_k(x)| > 1 \implies dr_g(x) = * \end{cases}$$

where $|\cdot|$ returns the cardinality of its argument and $*$ stands for a dummy category.

In words, dr_g returns either the index of the component function whose value is the highest, or the dummy category $*$ in case of ex æquo. In the case when the $g_k(x)$ are class posterior probability estimates, then dr is simply Bayes' estimated decision rule [31]. The qualifier *margin* refers to the fact that the generalization capabilities of such classifiers can be characterized by means of the values taken by the differences of the corresponding component functions. The use of the dummy category to avoid breaking ties is not central to the theory. Its main advantage rests in the fact that it keeps the reasoning and formulas as simple as possible.

With this definition at hand, the aim of the *learning process* is to minimize over \mathcal{G} the probability of error $P(dr_g(X) \neq Y)$. This probability can be reformulated in a handy way thanks to the introduction of additional functions.

Definition 4 (Class of functions $\mathcal{F}_{\mathcal{G}}$) Let \mathcal{G} be a class of functions satisfying Definition 3. For all $g \in \mathcal{G}$, the function f_g from $\mathcal{X} \times \llbracket 1, C \rrbracket$ into $[-M_{\mathcal{G}}, M_{\mathcal{G}}]$ is defined by:

$$\forall (x, k) \in \mathcal{X} \times \llbracket 1, C \rrbracket, \quad f_g(x, k) = \frac{1}{2} \left(g_k(x) - \max_{l \neq k} g_l(x) \right).$$

Then, the class $\mathcal{F}_{\mathcal{G}}$ is defined as follows:

$$\mathcal{F}_{\mathcal{G}} = \{f_g : g \in \mathcal{G}\}.$$

Definition 5 (Expected risk L) Let \mathcal{G} be a class of functions satisfying Definition 3 and let ϕ be the standard indicator loss function given by:

$$\forall t \in \mathbb{R}, \quad \phi(t) = \mathbb{1}_{\{t \leq 0\}}.$$

The expected risk of any function $g \in \mathcal{G}$, $L(g)$, is given by:

$$L(g) = \mathbb{E}_{(X,Y) \sim P} [\phi \circ f_g(X, Y)] = P(dr_g(X) \neq Y).$$

Its empirical risk measured on the m -sample \mathbf{Z}_m is:

$$L_m(g) = \mathbb{E}_{Z' \sim P_m} [\phi \circ f_g(Z')] = \frac{1}{m} \sum_{i=1}^m \phi \circ f_g(Z_i).$$

In order to take benefit from the fact that the classifiers of interest are margin ones, the sample-based estimate of performance which is actually used (involved in the different guaranteed risks) is obtained by substituting to ϕ a (dominating) margin loss/cost function. In this study, the definition used for those functions is the following one.

Definition 6 (Margin loss functions) A class of margin loss functions ϕ_{γ} parameterized by $\gamma \in (0, 1]$ is a class of nonincreasing functions from \mathbb{R} into $[0, 1]$ satisfying:

1. $\forall \gamma \in (0, 1], \quad \phi_{\gamma}(0) = 1 \wedge \phi_{\gamma}(\gamma) = 0;$
2. $\forall (\gamma, \gamma') \in (0, 1]^2, \quad \gamma < \gamma' \implies \forall t \in (0, \gamma), \quad \phi_{\gamma}(t) \leq \phi_{\gamma'}(t).$

Remark 1 The qualifier dominating is appropriate since we have for all $(\gamma, t) \in (0, 1] \times \mathbb{R}$, $\phi_{\gamma}(t) \geq \phi(t)$. The second property is especially useful to derive guaranteed risks holding uniformly for all values of γ . This can be achieved by means of Proposition 8 in [4]. It is noteworthy that these losses are not convex. They can even be discontinuous (whereas the definition used by Koltchinskii and Panchenko in [21] (Section 2) includes the Lipschitz property).

A risk obtained by substituting to ϕ a function ϕ_γ is named a margin risk.

Definition 7 (Margin risk L_γ) *Let \mathcal{G} be a class of functions satisfying Definition 3. For every (ordered) pair $(g, \gamma) \in \mathcal{G} \times (0, 1]$, the risk with margin γ of g , $L_\gamma(g)$, is defined as:*

$$L_\gamma(g) = \mathbb{E}_{Z \sim P} [\phi_\gamma \circ f_g(Z)].$$

$L_{\gamma,m}(g)$ designates the corresponding empirical risk, measured on the m -sample \mathbf{Z}_m :

$$L_{\gamma,m}(g) = \mathbb{E}_{Z' \sim P_m} [\phi_\gamma \circ f_g(Z')] = \frac{1}{m} \sum_{i=1}^m \phi_\gamma \circ f_g(Z_i).$$

Taking our inspiration from [4], we use margin loss functions in combination with a piecewise-linear squashing function. In short, the idea is to restrict the available information to what is relevant for the assessment of the prediction accuracy (the value of the margin loss is not affected), so as to optimize the way the introduction of the margin parameter γ is taken into account.

Definition 8 (Piecewise-linear squashing function π_γ) *For $\gamma \in (0, 1]$, the piecewise-linear squashing function π_γ is defined by:*

$$\forall t \in \mathbb{R}, \quad \pi_\gamma(t) = t \mathbb{1}_{\{t \in (0, \gamma]\}} + \gamma \mathbb{1}_{\{t > \gamma\}}.$$

This definition actually satisfies the aforementioned specification since we have:

$$\forall \gamma \in (0, 1], \quad \phi_\gamma \circ \pi_\gamma = \phi_\gamma.$$

Definition 9 (Class of functions $\mathcal{F}_{\mathcal{G},\gamma}$) *Let \mathcal{G} be a class of functions satisfying Definition 3 and $\mathcal{F}_{\mathcal{G}}$ the class of functions deduced from \mathcal{G} according to Definition 4. For every pair $(g, \gamma) \in \mathcal{G} \times (0, 1]$, the function $f_{g,\gamma}$ from $\mathcal{X} \times [1, C]$ into $[0, \gamma]$ is defined by:*

$$f_{g,\gamma} = \pi_\gamma \circ f_g.$$

Then, the class $\mathcal{F}_{\mathcal{G},\gamma}$ is defined as follows:

$$\mathcal{F}_{\mathcal{G},\gamma} = \{f_{g,\gamma} : g \in \mathcal{G}\}.$$

2.2 Scale-sensitive capacity measures

The guaranteed risks are ordinarily obtained in several main steps, corresponding to a basic supremum inequality and successive upper bounds on the capacity measure it involves,

each of which corresponds to a change of capacity measure. Although the measures which are central to this study are covering numbers, we start by giving the definition of the Rademacher complexity since it is the measure appearing first in the case of the L_2 -norm. For $n \in \mathbb{N}^*$, a Rademacher sequence σ_n is a sequence $(\sigma_i)_{1 \leq i \leq n}$ of independent random signs, i.e., independent and identically distributed random variables taking the values -1 and 1 with probability $\frac{1}{2}$ (symmetric Bernoulli or Rademacher random variables).

Definition 10 (Rademacher complexity) *Let $(\mathcal{T}, \mathcal{A}_{\mathcal{T}})$ be a measurable space and let T be a random variable with values in \mathcal{T} , distributed according to a probability measure P_T on $(\mathcal{T}, \mathcal{A}_{\mathcal{T}})$. For $n \in \mathbb{N}^*$, let $\mathbf{T}_n = (T_i)_{1 \leq i \leq n}$ be an n -sample made up of independent copies of T and let $\sigma_n = (\sigma_i)_{1 \leq i \leq n}$ be a Rademacher sequence. Let \mathcal{F} be a class of real-valued functions with domain \mathcal{T} . The empirical Rademacher complexity of \mathcal{F} is*

$$\hat{R}_n(\mathcal{F}) = \mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(T_i) \mid \mathbf{T}_n \right].$$

The Rademacher complexity of \mathcal{F} is

$$R_n(\mathcal{F}) = \mathbb{E}_{\mathbf{T}_n} [\hat{R}_n(\mathcal{F})] = \mathbb{E}_{\mathbf{T}_n \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(T_i) \right].$$

Remark 2 *The fact that the functional classes \mathcal{F} of interest can be uncountable calls for a specification. We make use of the standard convention (see for instance Formula (0.2) in [32]). Let $(T_s)_{s \in \mathcal{S}}$ be a stochastic process. Then,*

$$\mathbb{E} \left[\sup_{s \in \mathcal{S}} T_s \right] = \sup_{\{\bar{\mathcal{S}} \subset \mathcal{S}: |\bar{\mathcal{S}}| < +\infty\}} \mathbb{E} \left[\max_{s \in \bar{\mathcal{S}}} T_s \right].$$

The concept of covering number (ϵ -entropy), as well as the underlying concepts of ϵ -cover and ϵ -net, can be traced back to [20].

Definition 11 (ϵ -cover, ϵ -net, covering numbers, and ϵ -entropy) *Let (E, ρ) be a pseudo-metric space, $E' \subset E$ and $\epsilon \in \mathbb{R}_+^*$. An ϵ -cover of E' is a coverage of E' with open balls of radius ϵ the centers of which belong to E . These centers form an ϵ -net of E' . A proper ϵ -net of E' is an ϵ -net of E' included in E' . If E' has an ϵ -net of finite cardinality, then its covering number $\mathcal{N}(\epsilon, E', \rho)$ is the smallest cardinality of its ϵ -nets:*

$$\mathcal{N}(\epsilon, E', \rho) = \min \left\{ |E''| : (E'' \subset E) \wedge (\forall e \in E', \rho(e, E'') < \epsilon) \right\}.$$

If there is no such finite net, then the covering number is defined to be infinite. The corresponding logarithm, $\log_2(\mathcal{N}(\epsilon, E', \rho))$, is called the minimal ϵ -entropy of E' , or simply

the ϵ -entropy of E' . $\mathcal{N}^{(p)}(\epsilon, E', \rho)$ will designate a covering number of E' obtained by considering proper ϵ -nets only. In the finite case, we have thus:

$$\mathcal{N}^{(p)}(\epsilon, E', \rho) = \min \{|E''| : (E'' \subset E') \wedge (\forall e \in E', \rho(e, E'') < \epsilon)\}.$$

There is a close connection between covering and packing properties of bounded subsets in pseudo-metric spaces.

Definition 12 (ϵ -separation and packing numbers [20]) Let (E, ρ) be a pseudo-metric space and $\epsilon \in \mathbb{R}_+^*$. A set $E' \subset E$ is ϵ -separated if, for any distinct points e and e' in E' , $\rho(e, e') \geq \epsilon$. The ϵ -packing number of $E'' \subset E$, $\mathcal{M}(\epsilon, E'', \rho)$, is the maximal cardinality of an ϵ -separated subset of E'' , if such maximum exists. Otherwise, the ϵ -packing number of E'' is defined to be infinite.

In this study, the functional classes met are endowed with empirical (pseudo-)metrics derived from the L_p -norm.

Definition 13 (Pseudo-distance d_{p, \mathbf{t}_n}) Let \mathcal{F} be a class of real-valued functions on \mathcal{T} . For $n \in \mathbb{N}^*$, let $\mathbf{t}_n = (t_i)_{1 \leq i \leq n} \in \mathcal{T}^n$. Then,

$$\forall p \in \mathbb{N}^*, \forall (f, f') \in \mathcal{F}^2, d_{p, \mathbf{t}_n}(f, f') = \|f - f'\|_{L_p(\mu_{\mathbf{t}_n})} = \left(\frac{1}{n} \sum_{i=1}^n |f(t_i) - f'(t_i)|^p \right)^{\frac{1}{p}}$$

and

$$\forall (f, f') \in \mathcal{F}^2, d_{\infty, \mathbf{t}_n}(f, f') = \|f - f'\|_{L_\infty(\mu_{\mathbf{t}_n})} = \max_{1 \leq i \leq n} |f(t_i) - f'(t_i)|,$$

where $\mu_{\mathbf{t}_n}$ denotes the uniform (counting) probability measure on $\{t_i : 1 \leq i \leq n\}$.

Definition 14 (Uniform covering numbers [35] and uniform packing numbers [4])

Let \mathcal{F} be a class of real-valued functions on \mathcal{T} and $\bar{\mathcal{F}} \subset \mathcal{F}$. For $p \in \mathbb{N}^* \cup \{+\infty\}$, $\epsilon \in \mathbb{R}_+^*$, and $n \in \mathbb{N}^*$, the uniform covering number $\mathcal{N}_p(\epsilon, \bar{\mathcal{F}}, n)$ and the uniform packing number $\mathcal{M}_p(\epsilon, \bar{\mathcal{F}}, n)$ are defined as follows:

$$\begin{cases} \mathcal{N}_p(\epsilon, \bar{\mathcal{F}}, n) = \sup_{\mathbf{t}_n \in \mathcal{T}^n} \mathcal{N}(\epsilon, \bar{\mathcal{F}}, d_{p, \mathbf{t}_n}) \\ \mathcal{M}_p(\epsilon, \bar{\mathcal{F}}, n) = \sup_{\mathbf{t}_n \in \mathcal{T}^n} \mathcal{M}(\epsilon, \bar{\mathcal{F}}, d_{p, \mathbf{t}_n}) \end{cases}.$$

We define accordingly $\mathcal{N}_p^{(p)}(\epsilon, \bar{\mathcal{F}}, n)$ as:

$$\mathcal{N}_p^{(p)}(\epsilon, \bar{\mathcal{F}}, n) = \sup_{\mathbf{t}_n \in \mathcal{T}^n} \mathcal{N}^{(p)}(\epsilon, \bar{\mathcal{F}}, d_{p, \mathbf{t}_n}).$$

Our Sauer-Shelah lemma relates covering/packing numbers to a scale-sensitive generalization of the Vapnik-Chervonenkis (VC) dimension [34]: the fat-shattering dimension [18] also known as the γ -dimension.

Definition 15 (Fat-shattering dimension [18]) *Let \mathcal{F} be a class of functions from \mathcal{T} into $\mathcal{S} \subset \mathbb{R}$. For $\gamma \in \mathbb{R}_+^*$, a subset $s_{\mathcal{T}^n} = \{t_i : 1 \leq i \leq n\}$ of \mathcal{T} is said to be γ -shattered by \mathcal{F} if there is a vector $\mathbf{b}_n = (b_i)_{1 \leq i \leq n} \in \mathcal{S}^n$ such that, for every vector $\mathbf{l}_n = (l_i)_{1 \leq i \leq n} \in \{-1, 1\}^n$, there is a function $f_{\mathbf{l}_n} \in \mathcal{F}$ satisfying*

$$\forall i \in [1, n], \quad l_i (f_{\mathbf{l}_n}(t_i) - b_i) \geq \gamma.$$

The vector \mathbf{b}_n is called a witness to the γ -shattering. The fat-shattering dimension with margin γ of the class \mathcal{F} , $\gamma\text{-dim}(\mathcal{F})$, is the maximal cardinality of a subset of \mathcal{T} γ -shattered by \mathcal{F} , if such maximum exists. Otherwise, \mathcal{F} is said to have infinite fat-shattering dimension with margin γ .

Remark 3 *With the introduction of the set \mathcal{S} (and the constraint $\mathbf{b}_n \in \mathcal{S}^n$) in Definition 15, there is no need to make use of the strong dimension (Definition 3.1 in [1]). A difference with the definition used in [29] regards the concept of shattering. As most of the authors (see for instance [1]), we do not adopt the convention consisting in considering that the empty set can be shattered. Using the terminology of Mendelson and Vershynin (see Section 2.2 in [29]), the trivial center is not involved in our computations.*

Each of the generalized Sauer-Shelah lemmas in the literature is based on a main combinatorial result that involves a class of functions whose domain and codomain are finite sets. The first property is simply obtained by application of a restriction of the domain to the data at hand. As for the finiteness of the codomain, if needed, it is obtained by application of a discretization operator. The present study makes use of the following one, already employed, for instance, in [6].

Definition 16 (η -discretization operator) *Let \mathcal{F} be a class of functions from \mathcal{T} into $[-M_{\mathcal{F}}, M_{\mathcal{F}}]$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$. For $\eta \in \mathbb{R}_+^*$, define the η -discretization as an operator on \mathcal{F} such that:*

$$\begin{aligned} (\cdot)^{(\eta)} : \mathcal{F} &\longrightarrow \mathcal{F}^{(\eta)} \\ f &\mapsto f^{(\eta)} \\ \forall t \in \mathcal{T}, \quad f^{(\eta)}(t) &= \eta \left\lfloor \frac{f(t) + M_{\mathcal{F}}}{\eta} \right\rfloor \end{aligned}$$

where the floor function $\lfloor \cdot \rfloor$ is defined by:

$$\forall u \in \mathbb{R}, \lfloor u \rfloor = \max \{j \in \mathbb{Z} : j \leq u\}.$$

The finiteness of all the capacity measures considered in the sequel is ensured. Precisely, Theorem 2.5 in [1] (see also Theorem 2.4 in [28]) tells us that the fat-shattering dimension of a GC class is finite for every positive value of γ , and a corollary of the generalized Sauer-Shelah lemma is that the finiteness of this dimension implies the total boundedness.

3 L_p -norm Sauer-Shelah Lemma

Our master lemma is made up of two partial results. The first one, the *decomposition lemma*, relates the covering numbers of $\mathcal{F}_{\mathcal{G},\gamma}$ to those of the classes of component functions \mathcal{G}_k . The second one is the actual generalized Sauer-Shelah lemma.

3.1 Master lemma

Lemma 1 (Decomposition lemma) *Let \mathcal{G} be a class of functions satisfying Definition 3 and $\mathcal{F}_{\mathcal{G}}$ the class of functions deduced from \mathcal{G} according to Definition 4. For $\gamma \in (0, 1]$, let $\mathcal{F}_{\mathcal{G},\gamma}$ be the class of functions deduced from \mathcal{G} according to Definition 9. Then, for $\epsilon \in \mathbb{R}_+^*$, $m \in \mathbb{N}^*$, and $\mathbf{z}_m = ((x_i, y_i))_{1 \leq i \leq m} = (z_i)_{1 \leq i \leq m}$,*

$$\forall p \in \mathbb{N}^* \cup \{+\infty\}, \mathcal{N}^{(p)}(\epsilon, \mathcal{F}_{\mathcal{G},\gamma}, d_{p,\mathbf{z}_m}) \leq \mathcal{N}^{(p)}(\epsilon, \mathcal{F}_{\mathcal{G}}, d_{p,\mathbf{z}_m}) \leq \prod_{k=1}^C \mathcal{N}^{(p)}\left(\frac{\epsilon}{C^{\frac{1}{p}}}, \mathcal{G}_k, d_{p,\mathbf{x}_m}\right), \quad (1)$$

where $\mathbf{x}_m = (x_i)_{1 \leq i \leq m}$.

Proof The left-hand side inequality in Formula (1) is trivially true for $\epsilon > \gamma$. Otherwise, it is a direct consequence of the 1-Lipschitz property of the function π_γ . Similarly, the proof of the right-hand side inequality is nontrivial only for $\epsilon \leq 2M_{\mathcal{G}}$. We first derive it for a finite value of p only. For every function $g = (g_k)_{1 \leq k \leq C} \in \mathcal{G}$ and every element $z = (x, y) \in \mathcal{X} \times \llbracket 1, C \rrbracket$, let $k(g, z) \in \llbracket 1, C \rrbracket \setminus \{y\}$ be an index of category such that $f_g(z) = \frac{1}{2}(g_y(x) - g_{k(g,z)}(x))$. For all $k \in \llbracket 1, C \rrbracket$, let $\bar{\mathcal{G}}_k$ be a proper ϵ -net of \mathcal{G}_k with respect to the pseudo-metric d_{p,\mathbf{x}_m} such that $\bar{\mathcal{G}}_k$ is of cardinality $\mathcal{N}^{(p)}(\epsilon, \mathcal{G}_k, d_{p,\mathbf{x}_m})$. By construction, the cardinality of the class of functions $\bar{\mathcal{G}} = \prod_{k=1}^C \bar{\mathcal{G}}_k$ is $\prod_{k=1}^C \mathcal{N}^{(p)}(\epsilon, \mathcal{G}_k, d_{p,\mathbf{x}_m})$, and for every function $g = (g_k)_{1 \leq k \leq C} \in \mathcal{G}$, there exists a function $\bar{g} = (\bar{g}_k)_{1 \leq k \leq C} \in \bar{\mathcal{G}}$ such that:

$$\forall k \in \llbracket 1, C \rrbracket, d_{p,\mathbf{x}_m}(g_k, \bar{g}_k) < \epsilon. \quad (2)$$

By definition of the empirical pseudo-metric, for every $k \in \llbracket 1, C \rrbracket$ and every function $g_k \in \mathcal{G}_k$,

$$\begin{aligned} d_{p, \mathbf{x}_m}(g_k, \bar{g}_k) < \epsilon &\iff \left(\frac{1}{m} \sum_{i=1}^m |g_k(x_i) - \bar{g}_k(x_i)|^p \right)^{\frac{1}{p}} < \epsilon \\ &\implies \forall i \in \llbracket 1, m \rrbracket, \quad |g_k(x_i) - \bar{g}_k(x_i)| < m^{\frac{1}{p}} \epsilon \\ &\implies (|g_k(x_i) - \bar{g}_k(x_i)|)_{1 \leq i \leq m} = m^{\frac{1}{p}} \epsilon (\theta_{ki})_{1 \leq i \leq m} \end{aligned} \quad (3)$$

where $(\theta_{ki})_{1 \leq i \leq m} \in [0, 1]^m$. Furthermore, if $g_{k(g, z_i)}(x_i) \geq \bar{g}_{k(\bar{g}, z_i)}(x_i)$, then

$$\begin{aligned} |g_{k(g, z_i)}(x_i) - \bar{g}_{k(\bar{g}, z_i)}(x_i)| &= g_{k(g, z_i)}(x_i) - \bar{g}_{k(\bar{g}, z_i)}(x_i) \\ &\leq g_{k(g, z_i)}(x_i) - \bar{g}_{k(g, z_i)}(x_i) \\ &\leq |g_{k(g, z_i)}(x_i) - \bar{g}_{k(g, z_i)}(x_i)| \\ &\leq \theta_{k(g, z_i)i} m^{\frac{1}{p}} \epsilon. \end{aligned}$$

Symmetrically, $g_{k(g, z_i)}(x_i) \leq \bar{g}_{k(\bar{g}, z_i)}(x_i)$ implies that $|g_{k(g, z_i)}(x_i) - \bar{g}_{k(\bar{g}, z_i)}(x_i)| \leq \theta_{k(\bar{g}, z_i)i} m^{\frac{1}{p}} \epsilon$.

To sum up,

$$\forall i \in \llbracket 1, m \rrbracket, \quad |g_{k(g, z_i)}(x_i) - \bar{g}_{k(\bar{g}, z_i)}(x_i)| \leq \max(\theta_{k(g, z_i)i}, \theta_{k(\bar{g}, z_i)i}) m^{\frac{1}{p}} \epsilon. \quad (4)$$

For all $k \in \llbracket 1, C \rrbracket$, let $\boldsymbol{\theta}_k = (\theta_{ki})_{1 \leq i \leq m}$. Making use once more of (2) provides us with:

$$\forall k \in \llbracket 1, C \rrbracket, \quad \|\boldsymbol{\theta}_k\|_p < 1. \quad (5)$$

As a consequence,

$$\begin{aligned} d_{p, \mathbf{z}_m}(f_g, f_{\bar{g}}) &= \left(\frac{1}{m} \sum_{i=1}^m |f_g(z_i) - f_{\bar{g}}(z_i)|^p \right)^{\frac{1}{p}} \\ &= \frac{1}{2} \left(\frac{1}{m} \sum_{i=1}^m |g_{y_i}(x_i) - g_{k(g, z_i)}(x_i) - \bar{g}_{y_i}(x_i) + \bar{g}_{k(\bar{g}, z_i)}(x_i)|^p \right)^{\frac{1}{p}} \\ &\leq \frac{1}{2} \left(\frac{1}{m} \sum_{i=1}^m (|g_{y_i}(x_i) - \bar{g}_{y_i}(x_i)| + |g_{k(g, z_i)}(x_i) - \bar{g}_{k(\bar{g}, z_i)}(x_i)|)^p \right)^{\frac{1}{p}} \\ &\leq \frac{1}{2} \left(\sum_{i=1}^m (\theta_{y_i i} + \max(\theta_{k(g, z_i)i}, \theta_{k(\bar{g}, z_i)i}))^p \right)^{\frac{1}{p}} \epsilon \end{aligned} \quad (6)$$

$$\begin{aligned} &\leq \left(\sum_{i=1}^m \max_{1 \leq k \leq C} \theta_{ki}^p \right)^{\frac{1}{p}} \epsilon \\ &\leq \left(\sum_{k=1}^C \|\boldsymbol{\theta}_k\|_p^p \right)^{\frac{1}{p}} \epsilon \\ &< C^{\frac{1}{p}} \epsilon. \end{aligned} \quad (7)$$

Inequality (6) is obtained by application of (3) and (4), and Inequality (7) springs from Inequality (5). We have established that the set of functions $f_{\bar{g}}$ is a proper $(C^{\frac{1}{p}}\epsilon)$ -net of $\mathcal{F}_{\mathcal{G}}$ with respect to the pseudo-metric d_{p,\mathbf{z}_m} . Since its cardinality is at most that of $\bar{\mathcal{G}}$,

$$\forall \mathbf{z}_m \in (\mathcal{X} \times \llbracket 1, C \rrbracket)^m, \quad \mathcal{N}^{(p)}\left(C^{\frac{1}{p}}\epsilon, \mathcal{F}_{\mathcal{G}}, d_{p,\mathbf{z}_m}\right) \leq \prod_{k=1}^C \mathcal{N}^{(p)}(\epsilon, \mathcal{G}_k, d_{p,\mathbf{x}_m}).$$

The right-hand side inequality in Formula (1) then follows from performing a change of variable. The proof for the uniform convergence norm results from taking the limit when p goes to infinity. \blacksquare

The actual generalized Sauer-Shelah lemma is an extension of Lemma 3.5 in [1] and Lemma 8 in [6]. In the case when p is finite, then the upper bound is *dimension free* (does not depend on the number n of points) thanks to the implementation of the probabilistic extraction principle described in [29].

Lemma 2 (Generalized Sauer-Shelah lemma) *Let \mathcal{F} be a class of functions from \mathcal{T} into $[-M_{\mathcal{F}}, M_{\mathcal{F}}]$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$. \mathcal{F} is supposed to be a GC class. For $\epsilon \in (0, M_{\mathcal{F}}]$, let $d(\epsilon) = \epsilon \cdot \dim(\mathcal{F})$. Then for $\epsilon \in (0, 2M_{\mathcal{F}}]$ and $n \in \mathbb{N}^*$,*

$$\forall p \in \mathbb{N}^*, \quad \mathcal{M}_p(\epsilon, \mathcal{F}, n) \leq 2^{2(K_{\epsilon}(p)+1)} \left(\frac{6272eK_{\epsilon}(p)}{3} \left(\frac{2M_{\mathcal{F}}}{\epsilon} \right)^{2p+1} \right)^{2K_{\epsilon}(p)d(\frac{\epsilon}{45})}, \quad (8)$$

where $K_{\epsilon}(p) = \left\lceil (p+2) \log_2 \left(\left\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \right\rceil \right) \right\rceil$, and

$$\mathcal{M}_{\infty}(\epsilon, \mathcal{F}, n) \leq 2 \left(\frac{16M_{\mathcal{F}}^2 n}{\epsilon^2} \right)^{d(\frac{\epsilon}{4}) \log_2 \left(\frac{4M_{\mathcal{F}} n}{d(\frac{\epsilon}{4})\epsilon} \right)}. \quad (9)$$

Proof Since (9) is simply an instance of Lemma 3.5 in [1], we only prove (8). By definition,

$$\forall \mathbf{t}_n = (t_i)_{1 \leq i \leq n} \in \mathcal{T}^n, \quad \mathcal{M}(\epsilon, \mathcal{F}, d_{p,\mathbf{t}_n}) = \mathcal{M}(\epsilon, \mathcal{F}|_{\mathbf{t}_n}, d_{p,\mathbf{t}_n}),$$

where $\mathcal{F}|_{\mathbf{t}_n}$ is the set of the restrictions to \mathbf{t}_n of the functions in \mathcal{F} . Let \mathcal{F}_{ϵ} be, among the subsets of $\mathcal{F}|_{\mathbf{t}_n}$ ϵ -separated with respect to the pseudo-metric d_{p,\mathbf{t}_n} , a set of maximal cardinality. By definition,

$$|\mathcal{F}_{\epsilon}| = \mathcal{M}(\epsilon, \mathcal{F}|_{\mathbf{t}_n}, d_{p,\mathbf{t}_n}) = \mathcal{M}(\epsilon, \mathcal{F}_{\epsilon}, d_{p,\mathbf{t}_n}).$$

At this level, two cases must be considered.

First case Suppose that $|\mathcal{F}_\epsilon| \leq \exp(K_e(p) n \epsilon^{2p})$ where K_e is the function of p defined in Lemma 6. In that case, Lemma 6 applies, and we can set r equal to the smallest admissible value, $\frac{\ln(|\mathcal{F}_\epsilon|)}{K_e(p) \epsilon^{2p}}$, where \ln is the Neperian (or natural) logarithm. Consequently, there exists a subvector \mathbf{t}_q of \mathbf{t}_n of size

$$q \leq \frac{\ln(|\mathcal{F}_\epsilon|)}{K_e(p) \epsilon^{2p}} \quad (10)$$

such that \mathcal{F}_ϵ is $\left(\left(\frac{1}{2}\right)^{\frac{p+1}{p}} \epsilon\right)$ -separated with respect to the pseudo-metric d_{p, \mathbf{t}_q} , and thus, since $\min_{p \in \mathbb{N}^*} \left(\frac{1}{2}\right)^{\frac{p+1}{p}} = \frac{1}{4}$, $\frac{\epsilon}{4}$ -separated with respect to the same pseudo-metric. As a consequence,

$$|\mathcal{F}_\epsilon| = \mathcal{M}\left(\frac{\epsilon}{4}, \mathcal{F}_\epsilon, d_{p, \mathbf{t}_q}\right) = \mathcal{M}\left(\frac{\epsilon}{4}, \mathcal{F}_\epsilon|_{\mathbf{t}_q}, d_{p, \mathbf{t}_q}\right) = \left|\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right|.$$

For $\eta \in (0, \frac{\epsilon}{4})$, let $\left(\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right)^{(\eta)}$ be the image of $\mathcal{F}_\epsilon|_{\mathbf{t}_q}$ by the discretization operator $(\cdot)^\eta$. Since $\left|\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right| = \mathcal{M}\left(\frac{\epsilon}{4}, \mathcal{F}_\epsilon|_{\mathbf{t}_q}, d_{p, \mathbf{t}_q}\right)$, by application of Lemma 7,

$$\begin{aligned} \mathcal{M}\left(\frac{\epsilon}{4}, \mathcal{F}_\epsilon|_{\mathbf{t}_q}, d_{p, \mathbf{t}_q}\right) &= \mathcal{M}\left(\frac{\left(\left(\frac{\epsilon}{4}\right)^p - (\eta)^p\right)^{\frac{1}{p}}}{2}, \left(\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right)^{(\eta)}, d_{p, \mathbf{t}_q}\right) \\ &\leq \mathcal{M}\left(\frac{\epsilon - 4\eta}{8}, \left(\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right)^{(\eta)}, d_{p, \mathbf{t}_q}\right), \end{aligned}$$

where the inequality stems from the fact that

$$\min_{p \in \mathbb{N}^*} \left(\left(\frac{\epsilon}{4}\right)^p - (\eta)^p\right)^{\frac{1}{p}} = \frac{\epsilon}{4} - \eta.$$

For $N \in \mathbb{N}$ satisfying $N > \frac{56M_{\mathcal{F}}}{\epsilon}$, let us set $\eta = \frac{2M_{\mathcal{F}}}{N}$. Since $\left(\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right)^{\left(\frac{2M_{\mathcal{F}}}{N}\right)}$ is a class of functions whose domain has cardinality q and whose codomain is $\left\{2M_{\mathcal{F}} \frac{j}{N} : 0 \leq j \leq N\right\}$, Lemma 9 provides us with

$$|\mathcal{F}_\epsilon| \leq 2^{(p+2) \log_2(N)+1} \left(\frac{e(N-1)q}{d_1}\right)^{(p+2) \log_2(N)d_1}$$

where $d_1 = \left(\frac{1}{16} \left(\epsilon - \frac{56M_{\mathcal{F}}}{N}\right)\right) - \dim\left(\left(\mathcal{F}_\epsilon|_{\mathbf{t}_q}\right)^{\left(\frac{2M_{\mathcal{F}}}{N}\right)}\right)$. Thus, making use of the upper bound on q provided by (10),

$$\begin{aligned} |\mathcal{F}_\epsilon| &\leq 2^{(p+2) \log_2(N)+1} \left(\frac{e(N-1) \ln(|\mathcal{F}_\epsilon|)}{K_e(p) \epsilon^{2p} d_1}\right)^{(p+2) \log_2(N)d_1} \\ &\leq 2^{(p+2) \log_2(N)+1} \left(\frac{\ln(|\mathcal{F}_\epsilon|)}{d_1}\right)^{(p+2) \log_2(N)d_1} \left(\frac{e(N-1)}{K_e(p) \epsilon^{2p}}\right)^{(p+2) \log_2(N)d_1}. \end{aligned} \quad (11)$$

For all $r \in \mathbb{N}^*$, let h_r be the function on $[1, +\infty)$ mapping u to $2^{r-1}r! u^{\frac{1}{2}} - \ln^r(u)$. The function h_1 is positive on its domain and for all $r \in \mathbb{N}^*$, $h_r(1) > 0$. Since for all $r \geq 2$, $h'_r(u) = \frac{r}{u} h_{r-1}(u)$, proceeding by induction, one establishes that all the functions h_r are positive on their domain. Furthermore, for all $r \in \mathbb{N}^*$, $2^{r-1}r! \leq r^r$. Consequently, setting $K_{N,p} = \lceil (p+2) \log_2(N) \rceil$, where the ceiling function $\lceil \cdot \rceil$ is defined by:

$$\forall u \in \mathbb{R}, \lceil u \rceil = \min \{j \in \mathbb{Z} : j \geq u\},$$

we obtain

$$\begin{aligned} \left(\frac{\ln(|\mathcal{F}_\epsilon|)}{d_1} \right)^{(p+2) \log_2(N) d_1} &= \left(\ln^{(p+2) \log_2(N)} \left(|\mathcal{F}_\epsilon|^{\frac{1}{d_1}} \right) \right)^{d_1} \\ &\leq \left(\ln^{K_{N,p}} \left(|\mathcal{F}_\epsilon|^{\frac{1}{d_1}} \right) \right)^{d_1} \\ &< K_{N,p}^{K_{N,p} d_1} |\mathcal{F}_\epsilon|^{\frac{1}{2}}. \end{aligned} \quad (12)$$

A substitution of the right-hand side of (12) into (11) gives

$$|\mathcal{F}_\epsilon| \leq 2^{2(K_{N,p}+1)} \left(\frac{e(N-1) K_{N,p}}{K_e(p) \epsilon^{2p}} \right)^{2K_{N,p} d_1}.$$

To bound from above d_1 , N can be set equal to $\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \rceil$. Then,

$$d_1 \leq \left(\frac{\epsilon}{32} \right) - \dim \left(\left(\mathcal{F}_\epsilon|_{\mathbf{t}_q} \right)^{\left(\frac{2M_{\mathcal{F}}}{\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \rceil} \right)} \right) \quad (13)$$

$$\begin{aligned} &\leq \left(\frac{\epsilon}{32} - \frac{M_{\mathcal{F}}}{\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \rceil} \right) - \dim \left(\mathcal{F}_\epsilon|_{\mathbf{t}_q} \right) \\ &\leq \left(\frac{\epsilon}{32} - \frac{\epsilon}{112} \right) - \dim \left(\mathcal{F}_\epsilon|_{\mathbf{t}_q} \right) \\ &\leq \left(\frac{\epsilon}{45} \right) - \dim(\mathcal{F}). \end{aligned} \quad (14)$$

This sequence of computations makes use three times of the fact that the fat-shattering dimension is a nonincreasing function of the margin parameter. The transition from (13) to (14) is provided by Lemma 8. As a consequence,

$$|\mathcal{F}_\epsilon| \leq 2^{2(\lceil (p+2) \log_2(\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \rceil) \rceil + 1)} \left(\frac{112eM_{\mathcal{F}} \lceil (p+2) \log_2(\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \rceil) \rceil}{K_e(p) \epsilon^{2p+1}} \right)^{2 \lceil (p+2) \log_2(\lceil \frac{112M_{\mathcal{F}}}{\epsilon} \rceil) \rceil d(\frac{\epsilon}{45})}.$$

A substitution into the right-hand side of the value of $K_e(p)$ produces for $|\mathcal{F}_\epsilon|$, i.e., $\mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n})$, the same upper bound as that announced for $\mathcal{M}_p(\epsilon, \mathcal{F}, n)$. Thus, to conclude the proof of (8) under the assumption that $|\mathcal{F}_\epsilon| \leq \exp(K_e(p) n \epsilon^{2p})$, it suffices to notice that this upper bound does not depend on \mathbf{t}_n (it is even dimension free).

Second case Suppose conversely that $|\mathcal{F}_\epsilon| > \exp(K_e(p) n \epsilon^{2p})$, i.e.,

$$n < \frac{\ln(|\mathcal{F}_\epsilon|)}{K_e(p) \epsilon^{2p}}. \quad (15)$$

By application of Lemma 7, for $\eta \in (0, \epsilon)$,

$$\begin{aligned} |\mathcal{F}_\epsilon| &= \mathcal{M}\left(\frac{(\epsilon^p - \eta^p)^{\frac{1}{p}}}{2}, (\mathcal{F}_\epsilon)^{(\eta)}, d_{p, \mathbf{t}_n}\right) \\ &\leq \mathcal{M}\left(\frac{\epsilon - \eta}{2}, (\mathcal{F}_\epsilon)^{(\eta)}, d_{p, \mathbf{t}_n}\right). \end{aligned}$$

For $N \in \mathbb{N}$ satisfying $N > \frac{14M_{\mathcal{F}}}{\epsilon}$, let us set $\eta = \frac{2M_{\mathcal{F}}}{N}$. Since $(\mathcal{F}_\epsilon)^{(\frac{2M_{\mathcal{F}}}{N})}$ is a class of functions whose domain has cardinality n and whose codomain is $\{2M_{\mathcal{F}} \frac{j}{N} : 0 \leq j \leq N\}$, Lemma 9 provides us with

$$|\mathcal{F}_\epsilon| \leq 2^{(p+2)\log_2(N)+1} \left(\frac{e(N-1)n}{d_2}\right)^{(p+2)\log_2(N)d_2}$$

where $d_2 = \left(\frac{1}{4}\left(\epsilon - \frac{14M_{\mathcal{F}}}{N}\right)\right) - \dim\left((\mathcal{F}_\epsilon)^{(\frac{2M_{\mathcal{F}}}{N})}\right)$. The substitution of the upper bound on n provided by (15) into this bound produces

$$|\mathcal{F}_\epsilon| \leq 2^{2(K_{N,p}+1)} \left(\frac{e(N-1)K_{N,p}}{K_e(p)\epsilon^{2p}}\right)^{2K_{N,p}d_2}. \quad (16)$$

To bound from above d_2 , N can be set equal to $\left\lceil \frac{28M_{\mathcal{F}}}{\epsilon} \right\rceil$. Then, the line of reasoning used for d_1 leads to

$$d_2 \leq \left(\frac{\epsilon}{12}\right) - \dim(\mathcal{F}).$$

By substitution into (16) of the value of N and this upper bound on d_2 , an upper bound on $|\mathcal{F}_\epsilon|$ is obtained which is smaller than that provided by Inequality (8). \blacksquare

3.2 Comparison with the state of the art

In order to limit the complexity of the formula corresponding to finite values of p (Inequality (8)), the constants have systematically been derived by considering the “worst” case: $p = 1$. This implies that better constants can be obtained by focusing on the value of p of interest. If the resulting gain is all the more important as this value is large, it is already noticeable for $p = 2$. The result that compares directly with Lemma 2 is Theorem 3.2 in [28]. As Inequality (8), the corresponding bound is dimension free. The main difference rests in the dependency on the fat-shattering dimension. Whereas Inequality (8)

corresponds to a growth rate of the ϵ -entropy with this dimension which is linear, Theorem 3.2 in [28] exhibits an additional logarithmic multiplicative factor. Focusing on results derived for a specific L_p -norm, the literature provides us with one example of generalized Sauer-Shelah lemma based on the L_1 -norm: Lemma 1 in [5] (whose basic combinatorial result is Lemma 8 in [6]). However, this result is not dimension free (the growth rate of the ϵ -entropy with n is logarithmic). As for the L_2 -norm, the state of the art is provided by Theorem 1 in [29]. Since its original formulation involves unspecified universal constants, to make comparison possible, it is given below with explicit constants.

Lemma 3 (After Theorem 1 in [29]) *Let \mathcal{F} be a class of functions from \mathcal{T} into $[-M_{\mathcal{F}}, M_{\mathcal{F}}]$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$. \mathcal{F} is supposed to be a GC class. For $\epsilon \in (0, M_{\mathcal{F}}]$, let $d(\epsilon) = \epsilon\text{-dim}(\mathcal{F})$. Then for $\epsilon \in (0, 2M_{\mathcal{F}}]$ and $n \in \mathbb{N}^*$,*

$$\mathcal{M}_2(\epsilon, \mathcal{F}, n) \leq \left(3584e \left(\frac{2M_{\mathcal{F}}}{\epsilon} \right)^5 \right)^{4d(\frac{\epsilon}{96})}. \quad (17)$$

With this formulation at hand, it appears that even without optimizing the constants of Inequality (8) for the case $p = 2$, none of the two bounds is uniformly better than the other. The choice between them should primarily be based on the behaviour of the fat-shattering dimensions of interest.

4 Bound based on the L_{∞} -norm

The L_{∞} -norm plays a central part in the theory of bounds. Indeed, one can consider that it is already at the core of the initial result of Vapnik and Chervonenkis [34]. Focusing on margin classifiers, it is the norm used in Bartlett's seminal article [4].

4.1 State of the art

To the best of our knowledge, the state-of-the-art result is precisely a multi-class extension of Bartlett's result: Theorem 40 in [16]. It makes use of the same margin loss functions, defined as follows.

Definition 17 (Margin loss functions $\phi_{\infty, \gamma}$) *For $\gamma \in (0, 1]$, the margin loss function $\phi_{\infty, \gamma}$ is defined by:*

$$\forall t \in \mathbb{R}, \quad \phi_{\infty, \gamma}(t) = \mathbb{1}_{\{t < \gamma\}}.$$

The basic supremum inequality is a multi-class extension of Lemma 4 in [4], with the first symmetrization being derived from the basic lemma of Section 4.5.1 in [33].

Theorem 1 (After Theorem 22 in [16]) *Let \mathcal{G} be a class of functions satisfying Definition 3. For $\gamma \in (0, 1]$, let $\mathcal{F}_{\mathcal{G}, \gamma}$ be the class of functions deduced from \mathcal{G} according to Definition 9. For a fixed $\gamma \in (0, 1]$ and a fixed $\delta \in (0, 1)$, with P^m -probability at least $1 - \delta$, uniformly for every function $g \in \mathcal{G}$,*

$$L(g) \leq L_{\gamma, m}(g) + \sqrt{\frac{2}{m} \left(\ln \left(\mathcal{N}_{\infty}^{(p)} \left(\frac{\gamma}{2}, \mathcal{F}_{\mathcal{G}, \gamma}, 2m \right) \right) + \ln \left(\frac{2}{\delta} \right) \right)} + \frac{1}{m}, \quad (18)$$

where the margin loss function defining the empirical margin risk is $\phi_{\infty, \gamma}$ (Definition 17).

The pathway leading from this inequality to Theorem 40 in [16] consists in relating the covering number of interest to a γ - Ψ -dimension (see Definition 28 in [16]) of a class of vector-valued functions. The dependency on C varies with the choice of this dimension. In the case of the dimension which is the easiest to bound from above (by application of the pigeonhole principle), the margin Natarajan dimension, it is superlinear.

4.2 Improved dependency on C

Instead of working with vector-valued functions as in [16], it is more efficient to handle separately the classes of component functions. Starting from Inequality (18) and applying in sequence Lemma 1 (for $p = \infty$), Lemma 5 and Lemma 2 (Lemma 3.5 in [1]) produces the master theorem in the uniform convergence norm.

Theorem 2 *Let \mathcal{G} be a class of functions satisfying Definition 3. For $\epsilon \in (0, M_{\mathcal{G}}]$, let $d(\epsilon) = \max_{1 \leq k \leq C} \epsilon\text{-dim}(\mathcal{G}_k)$. For a fixed $\gamma \in (0, 1]$ and a fixed $\delta \in (0, 1)$, with P^m -probability at least $1 - \delta$, uniformly for every function $g \in \mathcal{G}$,*

$$L(g) \leq L_{\gamma, m}(g) + \sqrt{\frac{2}{m} \left(3Cd \left(\frac{\gamma}{8} \right) \ln^2 \left(\frac{128M_{\mathcal{G}}^2 m}{\gamma^2} \right) + \ln \left(\frac{2}{\delta} \right) \right)} + \frac{1}{m}.$$

Proof The sketch of the proof has been given at the beginning of the subsection. The detail makes use of the fact that

$$\forall k \in \llbracket 1, C \rrbracket, \quad \log_2 \left(\frac{16M_{\mathcal{G}} \epsilon m}{\left(\frac{\gamma}{8} \right)\text{-dim}(\mathcal{G}_k) \gamma} \right) \leq \frac{1}{\ln(2)} \ln \left(\frac{128M_{\mathcal{G}}^2 m}{\gamma^2} \right).$$

Thus,

$$\begin{aligned}
L(g) &\leq L_{\gamma,m}(g) + \sqrt{\frac{2}{m} \left(\sum_{k=1}^C \ln \left(\mathcal{N}_{\infty}^{(p)} \left(\frac{\gamma}{2}, \mathcal{G}_k, 2m \right) \right) + \ln \left(\frac{2}{\delta} \right) \right)} + \frac{1}{m} \\
&\leq L_{\gamma,m}(g) + \sqrt{\frac{2}{m} \left(\frac{2}{\ln(2)} \ln \left(\frac{128M_{\mathcal{G}}^2 m}{\gamma^2} \right) \sum_{k=1}^C \left(\frac{\gamma}{8} \right)^{-\dim(\mathcal{G}_k)} \ln \left(\frac{16M_{\mathcal{G}} e m}{\left(\frac{\gamma}{8} \right)^{-\dim(\mathcal{G}_k)} \gamma} \right) + \ln \left(\frac{2}{\delta} \right) \right)} + \frac{1}{m} \\
&\leq L_{\gamma,m}(g) + \sqrt{\frac{2}{m} \left(3 \ln^2 \left(\frac{128M_{\mathcal{G}}^2 m}{\gamma^2} \right) \sum_{k=1}^C \left(\frac{\gamma}{8} \right)^{-\dim(\mathcal{G}_k)} + \ln \left(\frac{2}{\delta} \right) \right)} + \frac{1}{m} \\
&\leq L_{\gamma,m}(g) + \sqrt{\frac{2}{m} \left(3Cd \left(\frac{\gamma}{8} \right) \ln^2 \left(\frac{128M_{\mathcal{G}}^2 m}{\gamma^2} \right) + \ln \left(\frac{2}{\delta} \right) \right)} + \frac{1}{m}.
\end{aligned}$$

■

4.3 Discussion

Under the assumption that $d(\epsilon)$ does not depend on C , Theorem 2 provides a guaranteed risk whose control term varies with C and m as a $O \left(\ln(m) \sqrt{\frac{C}{m}} \right)$. To sum up, the new bound exhibits the convergence rate of Theorem 40 in [16], whereas its control term grows only as the square root of C . Note that Lemma 19 in [36], which provides a bound with the same growth, holds for kernel multi-category classification methods only. We now establish an improvement of this kind with the L_2 -norm.

5 Bound based on the L_2 -norm

As in the case of the uniform convergence norm, the state-of-the-art result provides us not only with an element of comparison, but also with a starting point for the derivation of our guaranteed risk.

5.1 State of the art

The sharpest bound in the L_2 -norm is Theorem 3 in [23]. The margin loss function involved in this result is a standard one, the parameterized truncated hinge loss (that satisfies both Definition 6 and the definition used by Koltchinskii and Panchenko in [21]).

Definition 18 (Parameterized truncated hinge loss $\phi_{2,\gamma}$, Definition 4.3 in [30]) For $\gamma \in (0, 1]$, the parameterized truncated hinge loss $\phi_{2,\gamma}$ is defined by:

$$\forall t \in \mathbb{R}, \quad \phi_{2,\gamma}(t) = \mathbb{1}_{\{t \leq 0\}} + \left(1 - \frac{t}{\gamma}\right) \mathbb{1}_{\{t \in (0, \gamma]\}}.$$

This guaranteed risk is built upon a basic supremum inequality which is a partial result in the proof of Theorem 8.1 in [30] (with $\mathcal{F}_{\mathcal{G}}$ replaced with $\mathcal{F}_{\mathcal{G},\gamma}$).

Theorem 3 (After Theorem 8.1 in [30]) Let \mathcal{G} be a class of functions satisfying Definition 3. For $\gamma \in (0, 1]$, let $\mathcal{F}_{\mathcal{G},\gamma}$ be the class of functions deduced from \mathcal{G} according to Definition 9. For a fixed $\gamma \in (0, 1]$ and a fixed $\delta \in (0, 1)$, with P^m -probability at least $1 - \delta$, uniformly for every function $g \in \mathcal{G}$,

$$L(g) \leq L_{\gamma,m}(g) + \frac{2}{\gamma} R_m(\mathcal{F}_{\mathcal{G},\gamma}) + \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{2m}}$$

where the margin loss function defining the empirical margin risk is the parameterized truncated hinge loss (Definition 18).

Theorem 3 in [23] stems from Theorem 3 by application of the following lemma.

Lemma 4 Let \mathcal{G} be a class of functions satisfying Definition 3. For $\gamma \in (0, 1]$, let $\mathcal{F}_{\mathcal{G},\gamma}$ be the class of functions deduced from \mathcal{G} according to Definition 9. Then

$$R_m(\mathcal{F}_{\mathcal{G},\gamma}) \leq C R_m\left(\bigcup_{k=1}^C \mathcal{G}_k\right). \quad (19)$$

Many margin classifiers, including neural networks and kernel machines, satisfy the additional property that all the classes of component functions are identical, so that the growth with C of the upper bound on $R_m(\mathcal{F}_{\mathcal{G},\gamma})$ provided by (19) is linear. Furthermore, if the classifier is specifically a kernel machine, then it is well known that by combining the reproducing property with the Cauchy-Schwarz inequality, it is possible to obtain an upper bound on the Rademacher complexity which is a $O\left(m^{-\frac{1}{2}}\right)$ (see for instance Lemma 22 in [7]). Thus, for kernel machines, the control term of Kuznetsov's bound is a $O\left(\frac{C}{\sqrt{m}}\right)$. Kernel machines (with bounded range) satisfy Definition 3. This is easy to establish thanks to the characterization of the GC classes provided by Theorem 2.5 in [1]. The finiteness of the γ -dimension of a linear separator in a reproducing kernel Hilbert space is a well-known result, which appears, for instance, as a consequence of Theorem 4.6 in [8]. To sum up, the state-of-the-art result is a guaranteed risk whose control term is at best a $O\left(\frac{C}{\sqrt{m}}\right)$, for a specific family of classifiers among those satisfying Definition 3.

5.2 Improved dependency on C

Several results are available to bound from above the expected suprema of empirical processes (see for instance Chapters 1, 2, and 6 of [26]). We resort to the standard approach, especially efficient in the case of Rademacher processes, the application of Dudley's chaining method [14].

Theorem 4 (Chained bound on the Rademacher complexity of $\mathcal{F}_{\mathcal{G},\gamma}$) *Let \mathcal{G} be a class of functions satisfying Definition 3. For $\gamma \in (0, 1]$, let $\mathcal{F}_{\mathcal{G},\gamma}$ be the class of functions deduced from \mathcal{G} according to Definition 9. For $\epsilon \in (0, M_{\mathcal{G}}]$, let $d(\epsilon) = \max_{1 \leq k \leq C} \epsilon\text{-dim}(\mathcal{G}_k)$. Let h be a positive and decreasing function on \mathbb{N} such that $h(0) \geq \gamma$ and $h(1) \leq 2M_{\mathcal{G}}\sqrt{C}$. Then for all $N \in \mathbb{N}^*$,*

$$R_m(\mathcal{F}_{\mathcal{G},\gamma}) \leq h(N) + 4\sqrt{\frac{5C}{m}} \sum_{j=1}^N (h(j) + h(j-1)) \sqrt{d\left(\frac{h(j)}{96\sqrt{C}}\right) \ln\left(\frac{14M_{\mathcal{G}}\sqrt{C}}{h(j)}\right)}. \quad (20)$$

Proof The initial part of the proof of Formula (20) is the application of Theorem 6. Note that $\text{diam}(\mathcal{F}_{\mathcal{G},\gamma}) \leq \gamma$, justifying the hypothesis on $h(0)$. An advantage of working with $\mathcal{F}_{\mathcal{G},\gamma}$ instead of $\mathcal{F}_{\mathcal{G}}$ (directly) has thus been highlighted. The end of the proof consists in applying in sequence Lemma 1 (for $p = 2$), Lemma 5 and Lemma 3 (with $\epsilon = \frac{h(j)}{\sqrt{C}}$ and $3584e$ bounded from above by 7^5). ■

Thanks to the choice $h(j) = 2^{-j}\sqrt{C}\gamma$, under the assumption that $d(\epsilon)$ does not depend on C , then Theorem 4 provides a guaranteed risk whose control term grows linearly with C , a dependency at least as good as that of Theorem 3 in [23]. The improvement announced results from substituting to the hypothesis of GC classes a slightly stronger one.

Hypothesis 1 *We consider classes of functions \mathcal{G} satisfying Definition 3 plus the fact that there exists a pair $(d_{\mathcal{G}}, K_{\mathcal{G}}) \in \mathbb{N}^* \times \mathbb{R}_+^*$ such that*

$$\forall \epsilon \in (0, M_{\mathcal{G}}], \max_{1 \leq k \leq C} \epsilon\text{-dim}(\mathcal{G}_k) \leq K_{\mathcal{G}}\epsilon^{-d_{\mathcal{G}}}. \quad (21)$$

If Hypothesis 1 is satisfied, then the classes \mathcal{G}_k are *universal Donsker classes* [28]. Theorem 4.6 in [8] tells us that it is the case, with $d_{\mathcal{G}} = 2$, if each of the classes \mathcal{G}_k corresponds to the class of functions computed by a support vector machine (SVM) [11]. As a consequence, this is the case (with $d_{\mathcal{G}} = 2$) if \mathcal{G} is the class of functions computed by a multi-class SVM [17, 24, 13].

Theorem 5 Let \mathcal{G} be a class of functions satisfying Hypothesis 1. For $\gamma \in (0, 1]$, let $\mathcal{F}_{\mathcal{G}, \gamma}$ be the class of functions deduced from \mathcal{G} according to Definition 9.

If $d_{\mathcal{G}} = 1$, then

$$R_m(\mathcal{F}_{\mathcal{G}, \gamma}) \leq 160 \sqrt{\frac{30K_{\mathcal{G}}\gamma}{m}} C^{\frac{3}{4}} \left[\sqrt{\frac{\ln(F(C))}{2}} + \sqrt{\frac{\pi}{8}} F(C) \left(1 - \operatorname{erf}\left(\sqrt{\ln(F(C))}\right)\right) \right], \quad (22)$$

where

$$F(C) = 2 \sqrt{\frac{14M_{\mathcal{G}}}{\gamma}} C^{\frac{1}{4}}$$

and erf stands for the error function, i.e., $\operatorname{erf}(t) = \frac{2}{\sqrt{\pi}} \int_0^t e^{-u^2} du$.

If $d_{\mathcal{G}} = 2$, then

$$R_m(\mathcal{F}_{\mathcal{G}, \gamma}) \leq \frac{\gamma C^{\frac{3}{4}}}{\sqrt{m}} + 1152 \sqrt{\frac{5K_{\mathcal{G}}}{m}} C \left\lceil \frac{1}{2} \log_2 \left(\frac{m}{C} \right) \right\rceil \sqrt{\ln \left(\frac{14M_{\mathcal{G}} \sqrt{m}}{\gamma C^{\frac{1}{4}}} \right)}.$$

At last, if $d_{\mathcal{G}} > 2$, then

$$R_m(\mathcal{F}_{\mathcal{G}, \gamma}) \leq \sqrt{C} \left(\gamma \left(\frac{C}{m} \right)^{\frac{1}{d_{\mathcal{G}}}} + 8 \cdot 96^{\frac{d_{\mathcal{G}}}{2}} \left(2^{\frac{2}{d_{\mathcal{G}}-2}} + 1 \right) \cdot \gamma^{1-\frac{d_{\mathcal{G}}}{2}} \sqrt{5K_{\mathcal{G}}} \left(\frac{C}{m} \right)^{\frac{1}{d_{\mathcal{G}}}} \sqrt{\ln \left(\frac{14M_{\mathcal{G}}}{\gamma} \left(\frac{m}{C} \right)^{\frac{1}{d_{\mathcal{G}}}} \right)} \right). \quad (23)$$

Proof A substitution of Inequality (21) into Inequality (20) provides:

$$R_m(\mathcal{F}_{\mathcal{G}, \gamma}) \leq h(N) + 4 \cdot 96^{\frac{d_{\mathcal{G}}}{2}} \sqrt{\frac{5K_{\mathcal{G}}}{m}} C^{\frac{d_{\mathcal{G}}+2}{4}} \sum_{j=1}^N \frac{h(j) + h(j-1)}{h(j)^{\frac{d_{\mathcal{G}}}{2}}} \sqrt{\ln \left(\frac{14M_{\mathcal{G}} \sqrt{C}}{h(j)} \right)}. \quad (24)$$

At this point, we distinguish three cases according to the value taken by $d_{\mathcal{G}}$.

First case: $d_{\mathcal{G}} = 1$ This case is the only one for which the entropy integral of Formula (29) exists. Setting for all $j \in \mathbb{N}$, $h(j) = \gamma \cdot 2^{-2j}$, we obtain

$$R_m(\mathcal{F}_{\mathcal{G}, \gamma}) \leq 160 \sqrt{\frac{30K_{\mathcal{G}}\gamma}{m}} C^{\frac{3}{4}} \int_0^{\frac{1}{2}} \sqrt{\ln \left(\frac{14M_{\mathcal{G}} \sqrt{C}}{\gamma \epsilon^2} \right)} d\epsilon. \quad (25)$$

The computation of the integral gives

$$\int_0^{\frac{1}{2}} \sqrt{\ln \left(\frac{14M_{\mathcal{G}} \sqrt{C}}{\gamma \epsilon^2} \right)} d\epsilon = \sqrt{\frac{\ln(F(C))}{2}} + \frac{F(C)}{\sqrt{2}} \frac{\sqrt{\pi}}{2} \left(1 - \operatorname{erf}\left(\sqrt{\ln(F(C))}\right)\right). \quad (26)$$

Inequality (22) then results from a substitution of the right-hand side of (26) into (25).

Second case: $d_G = 2$ It stems from (24) that

$$R_m(\mathcal{F}_{\mathcal{G},\gamma}) \leq h(N) + 384\sqrt{\frac{5K_G}{m}}C \sum_{j=1}^N \frac{h(j) + h(j-1)}{h(j)} \sqrt{\ln\left(\frac{14M_G\sqrt{C}}{h(j)}\right)}.$$

For $N = \lceil \frac{1}{2} \log_2\left(\frac{m}{C}\right) \rceil$, we set $h(j) = \gamma C^{\frac{3}{4}} m^{-\frac{1}{2}} 2^{-j+N}$. Note that these choices are feasible since $N \in \mathbb{N}^*$ due to $m > C$, $h(0) \geq \gamma C^{\frac{1}{4}} > \gamma$, and $h(0) < 2\gamma C^{\frac{1}{4}} < 2M_G\sqrt{C}$. Then,

$$\begin{aligned} R_m(\mathcal{F}_{\mathcal{G},\gamma}) &\leq \frac{\gamma C^{\frac{3}{4}}}{\sqrt{m}} + 1152\sqrt{\frac{5K_G}{m}}C \sum_{j=1}^N \sqrt{\ln\left(\frac{14M_G\sqrt{m} \cdot 2^{j-N}}{\gamma C^{\frac{1}{4}}}\right)} \\ &\leq \frac{\gamma C^{\frac{3}{4}}}{\sqrt{m}} + 1152\sqrt{\frac{5K_G}{m}}C \left\lceil \frac{1}{2} \log_2\left(\frac{m}{C}\right) \right\rceil \sqrt{\ln\left(\frac{14M_G\sqrt{m}}{\gamma C^{\frac{1}{4}}}\right)}. \end{aligned}$$

Third case: $d_G > 2$ For $N = \left\lceil \frac{d_G-2}{2d_G} \log_2\left(\frac{m}{C}\right) \right\rceil$, let us set $h(j) = \gamma C^{\frac{1}{2} + \frac{1}{d_G}} m^{-\frac{1}{d_G}} 2^{\frac{2}{d_G-2}(-j+N)}$. Obviously, the constraints on N and the function h are once more satisfied. By substitution into (24), we get:

$$R_m(\mathcal{F}_{\mathcal{G},\gamma}) \leq \sqrt{C} \left(\gamma \left(\frac{C}{m}\right)^{\frac{1}{d_G}} + 4 \cdot 96^{\frac{d_G}{2}} \cdot \gamma^{1-\frac{d_G}{2}} \sqrt{5K_G} \left(\frac{C}{m}\right)^{\frac{1}{d_G}} \sqrt{\ln\left(\frac{14M_G}{\gamma} \left(\frac{m}{C}\right)^{\frac{1}{d_G}}\right)} S_N \right) \quad (27)$$

with

$$S_N = \sum_{j=1}^N \frac{2^{\frac{2}{d_G-2}(-j+N)} + 2^{\frac{2}{d_G-2}(-j+1+N)}}{2^{\frac{d_G}{d_G-2}(-j+N)}}.$$

Now,

$$\begin{aligned} S_N &= \left(2^{\frac{2}{d_G-2}} + 1 \right) \sum_{j=1}^N 2^{j-N} \\ &< 2 \left(2^{\frac{2}{d_G-2}} + 1 \right). \end{aligned}$$

Inequality (23) results from a substitution of this upper bound on S_N into (27). ■

5.3 Discussion

The implementation of Dudley's chaining method under Hypothesis 1 highlights the *phase transition* already identified by Mendelson in [28] (see also [27]). Besides this well-known phenomenon regarding the convergence rate, a parallel one can be noticed regarding the dependency on C . Indeed, if this dependency is always sublinear, as announced, it varies

significantly between \sqrt{C} and C , as a function of the value of $d_{\mathcal{G}}$. Its asymptotic value is \sqrt{C} . It is noteworthy that the behaviours observed are highly sensitive to the choice of the function h . We have already noticed in the beginning of the section that setting $h(j) = 2^{-j}\sqrt{C}\gamma$ has for consequence that the dependency on C is uniformly linear. Another example is instructive. In the case $d_{\mathcal{G}} = 1$, choosing $h(j) = \gamma \cdot 2^{-j}$ leads to

$$R_m(\mathcal{F}_{\mathcal{G},\gamma}) \leq 96\sqrt{\frac{30K_{\mathcal{G}}}{m}}C \int_0^{\frac{\gamma}{2\sqrt{C}}} \sqrt{\frac{1}{\epsilon} \ln\left(\frac{14M_{\mathcal{G}}}{\epsilon}\right)} d\epsilon.$$

6 Conclusions and ongoing research

An L_p -norm Sauer-Shelah lemma dedicated to margin multi-category classifiers whose classes of component functions are uniform Glivenko-Cantelli classes has been established. Its use makes it possible to improve the dependency on the number C of categories of the state-of-the-art guaranteed risks based on the L_{∞} -norm and the L_2 -norm. In both cases, this dependency becomes sublinear. Furthermore, in the favourable cases, the confidence interval can grow with C as slowly as a $O(\sqrt{C})$.

Our current work consists in continuing the unification of the approaches used to derive the bounds with respect to the different L_p -norms. The aim is to make the comparison of the resulting guaranteed risks more straightforward, as a step towards the characterization of the intrinsic complexity of the computation of polytomies. We also look for improvements resulting from the use of new tools from the theory of empirical processes. In that respect, the recent developments of the implementation of the chaining method appear promising.

Our results have been established under minimal assumptions regarding the pattern classification problem, the classifier and the margin loss function. Our future work will consist in assessing the benefit that one can derive from this study under different assumptions, such as those made in [24].

Acknowledgements The author would like to thank R. Vershynin for his explanations on the proof of Theorem 1 in [29] and A. Kontorovich for bringing to his attention the bounds in [30]. Thanks are also due to F. Lauer and K. Musayeva for carefully reading this manuscript. This work was partly funded by a CNRS research grant.

A Basic results and technical lemmas

Our formulation of Dudley's metric entropy bound, tailored for our needs, generalizes that established in transcripts of Bartlett's lectures which can be found online (see also [3]). The integral inequality appears as an instance of Corollary 13.2 in [10].

Theorem 6 (Dudley's metric entropy bound) *Let \mathcal{F} be a class of bounded real-valued functions on \mathcal{T} . For $n \in \mathbb{N}^*$, let $\mathbf{t}_n = (t_i)_{1 \leq i \leq n} \in \mathcal{T}^n$ and let $\text{diam}(\mathcal{F}) = \sup_{(f, f') \in \mathcal{F}^2} \|f - f'\|_{L_2(\mu_{\mathbf{t}_n})}$ be the diameter of \mathcal{F} in the $L_2(\mu_{\mathbf{t}_n})$ seminorm. Let h be a positive and decreasing function on \mathbb{N} such that $h(0) \geq \text{diam}(\mathcal{F})$. Then for $N \in \mathbb{N}^*$,*

$$\hat{R}_n(\mathcal{F}) \leq h(N) + 2 \sum_{j=1}^N (h(j) + h(j-1)) \sqrt{\frac{\ln(\mathcal{N}^{(p)}(h(j), \mathcal{F}, d_{2, \mathbf{t}_n}))}{n}} \quad (28)$$

and

$$\hat{R}_n(\mathcal{F}) \leq 12 \int_0^{\frac{1}{2} \text{diam}(\mathcal{F})} \sqrt{\frac{\ln(\mathcal{N}^{(p)}(\epsilon, \mathcal{F}, d_{2, \mathbf{t}_n}))}{n}} d\epsilon. \quad (29)$$

Proof For $j \in \mathbb{N}^*$, let $\bar{\mathcal{F}}_j$ be a proper $h(j)$ -net of \mathcal{F} with respect to d_{2, \mathbf{t}_n} such that $|\bar{\mathcal{F}}_j| = \mathcal{N}^{(p)}(h(j), \mathcal{F}, d_{2, \mathbf{t}_n})$. We set $\bar{\mathcal{F}}_0 = \{\bar{f}_0\}$ where \bar{f}_0 is any function in \mathcal{F} . Note that since $h(0)$ can be equal to $\text{diam}(\mathcal{F})$, the construction of $\bar{\mathcal{F}}_0$ does not ensure that this set is a proper $h(0)$ -net of \mathcal{F} with respect to d_{2, \mathbf{t}_n} (the minimum cardinality of such a net can be superior or equal to 2). The Rademacher process underlying the Rademacher complexity is centered, i.e.,

$$\forall f \in \mathcal{F}, \quad \mathbb{E}_{\sigma_n} \left[\frac{1}{\sqrt{n}} \sum_{i=1}^n \sigma_i f(t_i) \right] = 0.$$

Thus,

$$\hat{R}_n(\mathcal{F}) = \mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(t_i) - \bar{f}_0(t_i)) \right].$$

For each $f \in \mathcal{F}$ and each $j \in \mathbb{N}^*$, choose $\bar{f}_j \in \bar{\mathcal{F}}_j$ such that $\|f - \bar{f}_j\|_{L_2(\mu_{\mathbf{t}_n})} < h(j)$. Notice that

$$f - \bar{f}_0 = f - \bar{f}_N + \sum_{j=1}^N (\bar{f}_j - \bar{f}_{j-1}).$$

As a consequence, making use of the sub-additivity of the supremum function provides us with:

$$\hat{R}_n(\mathcal{F}) \leq \mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(t_i) - \bar{f}_N(t_i)) \right] + \sum_{j=1}^N \mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (\bar{f}_j(t_i) - \bar{f}_{j-1}(t_i)) \right]. \quad (30)$$

To bound from above the first term of the right-hand side of (30), we can make use in sequence of the Cauchy-Schwarz inequality and the definition of h .

$$\begin{aligned}
\mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(t_i) - \bar{f}_N(t_i)) \right] &\leq \mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \left\{ \left(\frac{1}{n} \sum_{i=1}^n \sigma_i^2 \right)^{\frac{1}{2}} \left(\frac{1}{n} \sum_{i=1}^n (f(t_i) - \bar{f}_N(t_i))^2 \right)^{\frac{1}{2}} \right\} \right] \\
&\leq \sup_{f \in \mathcal{F}} \|f - \bar{f}_N\|_{L_2(\mu_{\mathbf{t}_n})} \mathbb{E}_{\sigma_n} \left[\left(\frac{1}{n} \sum_{i=1}^n \sigma_i^2 \right)^{\frac{1}{2}} \right] \\
&< h(N).
\end{aligned} \tag{31}$$

As for the second term of the right-hand side of (30), we make use of Massart's finite class lemma (Lemma 5.2 in [25]). This calls for the derivation of an upper bound on $\left\| \frac{1}{n} (\bar{f}_j(t_i) - \bar{f}_{j-1}(t_i))_{1 \leq i \leq n} \right\|_2 = \frac{1}{\sqrt{n}} \|\bar{f}_j - \bar{f}_{j-1}\|_{L_2(\mu_{\mathbf{t}_n})}$ for all $j \in \llbracket 1, N \rrbracket$. This upper bound is obtained by application of Minkowski's inequality:

$$\begin{aligned}
\|\bar{f}_j - \bar{f}_{j-1}\|_{L_2(\mu_{\mathbf{t}_n})} &= \|\bar{f}_j - f + f - \bar{f}_{j-1}\|_{L_2(\mu_{\mathbf{t}_n})} \\
&\leq \|\bar{f}_j - f\|_{L_2(\mu_{\mathbf{t}_n})} + \|f - \bar{f}_{j-1}\|_{L_2(\mu_{\mathbf{t}_n})} \\
&< h(j) + h(j-1).
\end{aligned} \tag{32}$$

We can check that (32) still holds for $j = 1$ since

$$\begin{cases} \|\bar{f}_1 - f\|_{L_2(\mu_{\mathbf{t}_n})} < h(1) \\ \|f - \bar{f}_0\|_{L_2(\mu_{\mathbf{t}_n})} \leq \text{diam}(\mathcal{F}) \leq h(0) \end{cases} \implies \|\bar{f}_1 - f\|_{L_2(\mu_{\mathbf{t}_n})} + \|f - \bar{f}_0\|_{L_2(\mu_{\mathbf{t}_n})} < h(1) + h(0).$$

Applying Lemma 5.2 in [25] with (32) gives:

$$\begin{aligned}
\forall j \in \llbracket 1, N \rrbracket, \quad \mathbb{E}_{\sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (\bar{f}_j(t_i) - \bar{f}_{j-1}(t_i)) \right] &\leq \frac{h(j) + h(j-1)}{\sqrt{n}} \sqrt{2 \ln(|\bar{\mathcal{F}}_j| |\bar{\mathcal{F}}_{j-1}|)} \\
&\leq 2(h(j) + h(j-1)) \sqrt{\frac{\ln(|\bar{\mathcal{F}}_j|)}{n}}.
\end{aligned} \tag{33}$$

The substitution of (31) and (33) into (30) produces (28). Furthermore, setting for all

$j \in \mathbb{N}$, $h(j) = 2^{-j} \cdot \text{diam}(\mathcal{F})$, gives:

$$\begin{aligned} \hat{R}_n(\mathcal{F}) &\leq \text{diam}(\mathcal{F}) \left(2^{-N} + 6 \sum_{j=1}^N 2^{-j} \sqrt{\frac{\ln(\mathcal{N}^{(p)}(2^{-j} \cdot \text{diam}(\mathcal{F}), \mathcal{F}, d_{2, \mathbf{t}_n}))}{n}} \right) \\ &\leq \text{diam}(\mathcal{F}) \left(2^{-N} + 12 \sum_{j=1}^N (2^{-j} - 2^{-(j+1)}) \sqrt{\frac{\ln(\mathcal{N}^{(p)}(2^{-j} \cdot \text{diam}(\mathcal{F}), \mathcal{F}, d_{2, \mathbf{t}_n}))}{n}} \right) \end{aligned} \quad (34)$$

$$\leq 2^{-N} \cdot \text{diam}(\mathcal{F}) + 12 \int_{\frac{1}{2^{N+1}} \cdot \text{diam}(\mathcal{F})}^{\frac{1}{2} \cdot \text{diam}(\mathcal{F})} \sqrt{\frac{\ln(\mathcal{N}^{(p)}(\epsilon, \mathcal{F}, d_{2, \mathbf{t}_n}))}{n}} d\epsilon. \quad (35)$$

Inequality (35) springs from Inequality (34) since a covering number is a nonincreasing function of ϵ (on the interval $[2^{-(j+1)} \cdot \text{diam}(\mathcal{F}), 2^{-j} \cdot \text{diam}(\mathcal{F})]$, $\mathcal{N}^{(p)}(2^{-j} \cdot \text{diam}(\mathcal{F}), \mathcal{F}, d_{2, \mathbf{t}_n}) \leq \mathcal{N}^{(p)}(\epsilon, \mathcal{F}, d_{2, \mathbf{t}_n})$). Inequality (29) is simply the asymptotic formulation of Inequality (35) (for N going to infinity). \blacksquare

Lemma 5 (After Theorem IV in [20]) *Let (E, ρ) be a pseudo-metric space. For every totally bounded set $E' \subset E$ and $\epsilon \in \mathbb{R}_+^*$,*

$$\mathcal{N}^{(p)}(\epsilon, E', \rho) \leq \mathcal{M}(\epsilon, E', \rho).$$

Lemma 6 (After Lemma 13 in [29]) *Let $\mathcal{T} = \{t_i : 1 \leq i \leq n\}$ be a finite set and $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$. Let \mathcal{F} be a finite class of functions from \mathcal{T} into $[-M_{\mathcal{F}}, M_{\mathcal{F}}]$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$. Let $p \in \mathbb{N}^*$. Assume that for some $\epsilon \in (0, 2M_{\mathcal{F}}]$, \mathcal{F} is ϵ -separated with respect to the pseudo-metric d_{p, \mathbf{t}_n} . If $r \in [1, n]$ is such that $|\mathcal{F}| \leq \exp(K_e(p) r \epsilon^{2p})$ with*

$$K_e(p) = \frac{3}{112(2M_{\mathcal{F}})^{2p}},$$

then there exists a subvector \mathbf{t}_q of \mathbf{t}_n of size $q \leq r$ such that \mathcal{F} is $\left(\left(\frac{1}{2}\right)^{\frac{p+1}{p}} \epsilon\right)$ -separated with respect to the pseudo-metric d_{p, \mathbf{t}_q} .

Proof Let us set $\mathcal{F} = \{f_j : 1 \leq j \leq |\mathcal{F}|\}$ and $\mathcal{D}_{\mathcal{F}} = \{f_j - f_{j'} : 1 \leq j < j' \leq |\mathcal{F}|\}$. The set $\mathcal{D}_{\mathcal{F}}$ has cardinality $|\mathcal{D}_{\mathcal{F}}| < \frac{1}{2} |\mathcal{F}|^2$. Fix $r \in [1, n]$ satisfying the assumptions of the lemma and let $(\epsilon_i)_{1 \leq i \leq n}$ be a sequence of n independent Bernoulli random variables with common expectation $\mu = \frac{r}{2n}$. Then, by application of the ϵ -separation property, for every

δ_f in $\mathcal{D}_{\mathcal{F}}$,

$$\begin{aligned} \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p < \mu \left(\frac{\epsilon}{2} \right)^p \right) &\leq \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n (\mu - \epsilon_i) |\delta_f(t_i)|^p > \left(1 - \frac{1}{2^p} \right) \mu \epsilon^p \right) \\ &\leq \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n (\mu - \epsilon_i) |\delta_f(t_i)|^p > \frac{1}{2} \mu \epsilon^p \right). \end{aligned} \quad (36)$$

Since by construction, for all $i \in \llbracket 1, n \rrbracket$, $\mathbb{E}[(\mu - \epsilon_i) |\delta_f(t_i)|^p] = 0$ and $|\mu - \epsilon_i| |\delta_f(t_i)|^p \leq (2M_{\mathcal{F}})^p (1 - \mu) < (2M_{\mathcal{F}})^p$ with probability one, the right-hand side of (36) can be bounded from above thanks to Bernstein's inequality [9]. Given that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[(\mu - \epsilon_i)^2 \delta_f(t_i)^{2p} \right] \leq (2M_{\mathcal{F}})^{2p} \mu (1 - \mu) < (2M_{\mathcal{F}})^{2p} \mu,$$

we obtain

$$\begin{aligned} \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p < \mu \left(\frac{\epsilon}{2} \right)^p \right) &\leq \exp \left(- \frac{3\mu n \epsilon^{2p}}{4 \left(6(2M_{\mathcal{F}})^{2p} + (2M_{\mathcal{F}})^p \epsilon^p \right)} \right) \\ &\leq \exp \left(- \frac{3r \epsilon^{2p}}{56 (2M_{\mathcal{F}})^{2p}} \right) \\ &\leq \exp(-2K_e(p) r \epsilon^{2p}). \end{aligned}$$

Therefore, given the assumption on r , applying the union bound provides us with:

$$\begin{aligned} \mathbb{P} \left(\exists \delta_f \in \mathcal{D}_{\mathcal{F}} : \left(\frac{1}{r} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p \right)^{\frac{1}{p}} < \left(\frac{1}{2} \right)^{\frac{p+1}{p}} \epsilon \right) &= \mathbb{P} \left(\exists \delta_f \in \mathcal{D}_{\mathcal{F}} : \frac{1}{n} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p < \mu \left(\frac{\epsilon}{2} \right)^p \right) \\ &\leq \sum_{\delta_f \in \mathcal{D}_{\mathcal{F}}} \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p < \mu \left(\frac{\epsilon}{2} \right)^p \right) \\ &\leq |\mathcal{D}_{\mathcal{F}}| \cdot \exp(-2K_e(p) r \epsilon^{2p}) \\ &< \frac{1}{2} \exp^2(K_e(p) r \epsilon^{2p}) \cdot \exp(-2K_e(p) r \epsilon^{2p}) \\ &< \frac{1}{2}. \end{aligned} \quad (37)$$

Moreover, if \mathcal{S}_1 is the random set $\{i \in \llbracket 1, n \rrbracket : \epsilon_i = 1\}$, then by Markov's inequality,

$$\mathbb{P}(|\mathcal{S}_1| > r) = \mathbb{P} \left(\sum_{i=1}^n \epsilon_i > r \right) \leq \frac{1}{2}. \quad (38)$$

Combining (37) and (38) by means of the union bound provides us with

$$\mathbb{P} \left\{ \left(\exists \delta_f \in \mathcal{D}_{\mathcal{F}} : \left(\frac{1}{r} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p \right)^{\frac{1}{p}} < \left(\frac{1}{2} \right)^{\frac{p+1}{p}} \epsilon \right) \vee (|\mathcal{S}_1| > r) \right\} < 1$$

or equivalently

$$\mathbb{P} \left\{ \left(\forall \delta_f \in \mathcal{D}_{\mathcal{F}} : \left(\frac{1}{r} \sum_{i=1}^n \epsilon_i |\delta_f(t_i)|^p \right)^{\frac{1}{p}} \geq \left(\frac{1}{2} \right)^{\frac{p+1}{p}} \epsilon \right) \wedge (|S_1| \leq r) \right\} > 0$$

which implies that

$$\mathbb{P} \left\{ \left(\forall \delta_f \in \mathcal{D}_{\mathcal{F}} : \|\delta_f\|_{L_p(\mu_{(t_i)_{i \in S_1}})} \geq \left(\frac{1}{2} \right)^{\frac{p+1}{p}} \epsilon \right) \wedge (|S_1| \leq r) \right\} > 0.$$

This translates into the fact that there exists a subvector \mathbf{t}_q of \mathbf{t}_n of size $q \leq r$ such that the class \mathcal{F} is $\left(\left(\frac{1}{2} \right)^{\frac{p+1}{p}} \epsilon \right)$ -separated with respect to the pseudo-metric d_{p, \mathbf{t}_q} , i.e., our claim. \blacksquare

Lemma 7 *Let \mathcal{F} be a class of functions from \mathcal{T} into $[-M_{\mathcal{F}}, M_{\mathcal{F}}]$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$. For $n \in \mathbb{N}^*$, let $\mathbf{t}_n = (t_i)_{1 \leq i \leq n} \in \mathcal{T}^n$. For all $\epsilon \in (0, 2M_{\mathcal{F}}]$, all $\eta \in (0, \epsilon)$, and all $p \in \mathbb{N}^*$, if a subset of \mathcal{F} is ϵ -separated in the pseudo-metric d_{p, \mathbf{t}_n} , then the η -discretization operator acts on it as an injective mapping and the image obtained is a set $\left(\frac{(\epsilon^p - \eta^p)^{\frac{1}{p}}}{2} \right)$ -separated in the same pseudo-metric. As a consequence,*

$$\forall p \in \mathbb{N}^*, \quad \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}) \leq \mathcal{M} \left(\frac{(\epsilon^p - \eta^p)^{\frac{1}{p}}}{2}, \mathcal{F}^{(\eta)}, d_{p, \mathbf{t}_n} \right).$$

Proof Proving Lemma 7 amounts to establishing that

$$\forall (f_1, f_2) \in \mathcal{F}^2, \quad (d_{p, \mathbf{t}_n}(f_1, f_2) \geq \epsilon) \wedge (\eta \in (0, \epsilon)) \implies d_{p, \mathbf{t}_n}(f_1^{(\eta)}, f_2^{(\eta)}) \geq \frac{(\epsilon^p - \eta^p)^{\frac{1}{p}}}{2}. \quad (39)$$

For $i \in \llbracket 1, n \rrbracket$, let $\delta_i = \left(f_1^{(\eta)}(t_i) - f_2^{(\eta)}(t_i) \right)$ and $\delta'_i = f_1(t_i) - f_2(t_i) - \delta_i$. By construction, there exists $N_i \in \mathbb{N}$ such that $|\delta_i| = \eta N_i$, and $|\delta'_i| < \eta$. If $N_i > 0$, then $|\delta_i| + |\delta'_i| < 2|\delta_i|$, otherwise $|\delta_i| + |\delta'_i| < \eta$, with the consequence that in all cases, $(|\delta_i| + |\delta'_i|)^p < (2|\delta_i|)^p + \eta^p$.

Thus,

$$\begin{aligned} (d_{p, \mathbf{t}_n}(f_1, f_2) \geq \epsilon) \wedge (\eta \in (0, \epsilon)) &\implies \frac{1}{n} \sum_{i=1}^n |\delta_i + \delta'_i|^p \geq \epsilon^p \\ &\implies \frac{1}{n} \sum_{i=1}^n (|\delta_i| + |\delta'_i|)^p \geq \epsilon^p \\ &\implies \frac{1}{n} \sum_{i=1}^n (2|\delta_i|)^p + \eta^p \geq \epsilon^p \\ &\implies \left(2d_{p, \mathbf{t}_n}(f_1^{(\eta)}, f_2^{(\eta)}) \right)^p + \eta^p \geq \epsilon^p \\ &\implies d_{p, \mathbf{t}_n}(f_1^{(\eta)}, f_2^{(\eta)}) \geq \frac{(\epsilon^p - \eta^p)^{\frac{1}{p}}}{2}. \end{aligned}$$

To sum up, we have established (39), i.e., the lemma. ■

Lemma 8 *Let \mathcal{F} be a class of functions from \mathcal{T} into $[-M_{\mathcal{F}}, M_{\mathcal{F}}]$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$. For all $\epsilon \in (0, M_{\mathcal{F}}]$ and all $\eta \in (0, 2\epsilon)$,*

$$\epsilon\text{-dim}\left(\mathcal{F}^{(\eta)}\right) \leq \left(\epsilon - \frac{\eta}{2}\right)\text{-dim}(\mathcal{F}).$$

Proof To prove Lemma 8, it suffices to notice that

$$\begin{aligned} f^{(\eta)}(t) - b \geq \epsilon &\implies \eta \left\lfloor \frac{f(t) + M_{\mathcal{F}}}{\eta} \right\rfloor - b \geq \epsilon \\ &\implies f(t) + M_{\mathcal{F}} - b \geq \epsilon \\ &\implies f(t) - \left(b + \frac{\eta}{2} - M_{\mathcal{F}}\right) \geq \epsilon - \frac{\eta}{2} \end{aligned}$$

and

$$\begin{aligned} f^{(\eta)}(t) - b \leq -\epsilon &\implies \eta \left\lfloor \frac{f(t) + M_{\mathcal{F}}}{\eta} \right\rfloor - b \leq -\epsilon \\ &\implies f(t) + M_{\mathcal{F}} - \eta - b \leq -\epsilon \\ &\implies f(t) - \left(b + \frac{\eta}{2} - M_{\mathcal{F}}\right) \leq -\left(\epsilon - \frac{\eta}{2}\right). \end{aligned}$$

■

In the framework of this study, the main combinatorial result evoqued in Section 2.2 is the following lemma, which extends Lemma 8 in [6].

Lemma 9 *Let $\mathcal{T} = \{t_i : 1 \leq i \leq n\}$ be a finite set and $\mathbf{t}_n = (t_i)_{1 \leq i \leq n}$. Let \mathcal{F} be a class of functions from \mathcal{T} into $\mathcal{S} = \left\{2M_{\mathcal{F}}\frac{j}{N} : 0 \leq j \leq N\right\}$ with $M_{\mathcal{F}} \in \mathbb{R}_+^*$ and $N \in \mathbb{N} \setminus \llbracket 0, 3 \rrbracket$. For $\epsilon \in \left(\frac{6M_{\mathcal{F}}}{N}, 2M_{\mathcal{F}}\right]$, let $d = \left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)\text{-dim}(\mathcal{F})$. Then*

$$\forall p \in \mathbb{N}^*, \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}) < 2^{(p+2)\log_2(N)+1} \left(\frac{e(N-1)n}{d}\right)^{(p+2)\log_2(N)d}. \quad (40)$$

Proof First, note that

$$\begin{cases} d \geq 1 \\ \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}) \geq 2 \end{cases} \implies \begin{cases} \epsilon \in \left(\frac{6M_{\mathcal{F}}}{N}, 2M_{\mathcal{F}} + \frac{6M_{\mathcal{F}}}{N}\right] \\ \epsilon \in (0, 2M_{\mathcal{F}}] \end{cases} \implies \begin{cases} \epsilon \in \left(\frac{6M_{\mathcal{F}}}{N}, 2M_{\mathcal{F}}\right] \\ N > 3 \end{cases}.$$

For $q \in \llbracket 1, d \rrbracket$, let the pair $(s_{\mathcal{T}^q}, \mathbf{b}_q)$ be such that $s_{\mathcal{T}^q}$ is a subset of \mathcal{T} of cardinality q and $\mathbf{b}_q \in (\mathcal{S} \setminus \{0, 2M_{\mathcal{F}}\})^q$. Such a pair will be said to be γ -shattered by a subset of

\mathcal{F} if $s_{\mathcal{T}^q}$ is γ -shattered by this subset and \mathbf{b}_q is a witness to this shattering. Setting $K = \sum_{j=0}^d \binom{n}{j} (N-1)^j$, the number of such pairs is equal to $\sum_{j=1}^d \binom{n}{j} (N-1)^j$, i.e., to $K-1$. Fix $\epsilon \in \left(\frac{6M_{\mathcal{F}}}{N}, 2M_{\mathcal{F}}\right]$ and $p \in \mathbb{N}^*$. For each $r \in \llbracket 2, \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}) \rrbracket$, let $\text{shat}(r)$ be the maximum integer such that any subset of \mathcal{F} of cardinality r which is ϵ -separated in the metric d_{p, \mathbf{t}_n} $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shatters at least $\text{shat}(r)$ pairs $(s_{\mathcal{T}^q}, \mathbf{b}_q)$. Obviously, the function shat is nondecreasing. We now establish that $\text{shat}(2) \geq 1$. Indeed, let $\{f_+, f_-\}$ be a subset of \mathcal{F} ϵ -separated in the metric d_{p, \mathbf{t}_n} . By definition,

$$\left(\frac{1}{n} \sum_{i=1}^n |f_+(t_i) - f_-(t_i)|^p\right)^{\frac{1}{p}} \geq \epsilon,$$

with the consequence that there exists $i_0 \in \llbracket 1, n \rrbracket$ such that $|f_+(t_{i_0}) - f_-(t_{i_0})| \geq \epsilon$. Without loss of generality, we make the hypothesis that $f_+(t_{i_0}) - f_-(t_{i_0}) \geq \epsilon$. Then,

$$\begin{aligned} f_+(t_{i_0}) - \frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (f_+(t_{i_0}) + f_-(t_{i_0})) \right\rfloor &\geq f_+(t_{i_0}) - \frac{1}{2} (f_+(t_{i_0}) + f_-(t_{i_0})) \\ &\geq \frac{1}{2} (f_+(t_{i_0}) - f_-(t_{i_0})) \\ &\geq \frac{\epsilon}{2} \\ &\geq \frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N} \end{aligned}$$

and

$$\begin{aligned} f_-(t_{i_0}) - \frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (f_+(t_{i_0}) + f_-(t_{i_0})) \right\rfloor &\leq f_-(t_{i_0}) - \frac{1}{2} (f_+(t_{i_0}) + f_-(t_{i_0})) + \frac{2M_{\mathcal{F}}}{N} \\ &\leq \frac{1}{2} (f_-(t_{i_0}) - f_+(t_{i_0})) + \frac{2M_{\mathcal{F}}}{N} \\ &\leq -\frac{\epsilon}{2} + \frac{2M_{\mathcal{F}}}{N} \\ &\leq -\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right). \end{aligned}$$

Since $\frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (f_+(t_{i_0}) + f_-(t_{i_0})) \right\rfloor \in \mathcal{S} \setminus \{0, 2M_{\mathcal{F}}\}$, we have established that the set $\{f_+, f_-\}$ $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shatters $(\{t_{i_0}\}, \mathbf{b}_1)$ with $\mathbf{b}_1 = \left(\frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (f_+(t_{i_0}) + f_-(t_{i_0})) \right\rfloor\right)$, which concludes the proof of $\text{shat}(2) \geq 1$. Choose an even $r \in \llbracket 2, \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}) \rrbracket$ and let $\bar{\mathcal{F}}$ be a subset of \mathcal{F} of cardinality r ϵ -separated in the metric d_{p, \mathbf{t}_n} . Split $\bar{\mathcal{F}}$ arbitrarily into $\frac{r}{2}$ pairs. For each such pair (f_+, f_-) , let

$$\text{ind}(f_+, f_-) = \left| \left\{ i \in \llbracket 1, n \rrbracket : |f_+(t_i) - f_-(t_i)| \geq \epsilon - \frac{2M_{\mathcal{F}}}{N} \right\} \right|.$$

Then

$$\begin{aligned}
d_{p, \mathbf{t}_n}(f_+, f_-) &= \left(\frac{1}{n} \sum_{i=1}^n |f_+(t_i) - f_-(t_i)|^p \right)^{\frac{1}{p}} \\
&\leq \left\{ \frac{1}{n} \left[\text{ind}(f_+, f_-) (2M_{\mathcal{F}})^p + (n - \text{ind}(f_+, f_-)) \left(\epsilon - \frac{2M_{\mathcal{F}}}{N} \right)^p \right] \right\}^{\frac{1}{p}} \\
&\leq \left[\frac{\text{ind}(f_+, f_-)}{n} (2M_{\mathcal{F}})^p + \left(\epsilon - \frac{2M_{\mathcal{F}}}{N} \right)^p \right]^{\frac{1}{p}}.
\end{aligned}$$

By hypothesis, $d_{p, \mathbf{t}_n}(f_+, f_-) \geq \epsilon > \frac{6M_{\mathcal{F}}}{N}$, which implies that

$$\begin{aligned}
\text{ind}(f_+, f_-) &\geq \frac{n}{(2M_{\mathcal{F}})^p} \left[\epsilon^p - \left(\epsilon - \frac{2M_{\mathcal{F}}}{N} \right)^p \right] \\
&\geq \frac{n}{N (2M_{\mathcal{F}})^{p-1}} \sum_{j=0}^{p-1} \epsilon^{p-j-1} \left(\epsilon - \frac{2M_{\mathcal{F}}}{N} \right)^j \\
&\geq \frac{n}{N (2M_{\mathcal{F}})^{p-1}} \sum_{j=0}^{p-1} \left(\frac{6M_{\mathcal{F}}}{N} \right)^{p-j-1} \left(\frac{4M_{\mathcal{F}}}{N} \right)^j \\
&\geq \frac{n}{N^p} \sum_{j=0}^{p-1} 3^{p-j-1} 2^j \\
&\geq \frac{3^p - 2^p}{N^p} n \\
&\geq \frac{n}{N^p}.
\end{aligned}$$

Thus, each pair (f_+, f_-) has at least $\frac{n}{N^p}$ indices i such that $|f_+(t_i) - f_-(t_i)| \geq \epsilon - \frac{2M_{\mathcal{F}}}{N}$. Applying the pigeonhole principle, there is at least one index i_0 such that at least $\lceil \frac{rn}{2N^p n} \rceil = \lceil \frac{r}{2N^p} \rceil$ pairs (f_+, f_-) satisfy $|f_+(t_{i_0}) - f_-(t_{i_0})| \geq \epsilon - \frac{2M_{\mathcal{F}}}{N}$. Keeping in mind that $\epsilon - \frac{2M_{\mathcal{F}}}{N} > \frac{4M_{\mathcal{F}}}{N}$, it is easy to establish that there are less than $\frac{N^2}{2}$ different pairs $(u_1, u_2) \in \mathcal{S}^2$ such that $|u_1 - u_2| \geq \epsilon - \frac{2M_{\mathcal{F}}}{N}$. Thus, applying once more the pigeonhole principle, there are at least $\lceil \frac{r}{N^{p+2}} \rceil$ pairs (f_+, f_-) such that $|f_+(t_{i_0}) - f_-(t_{i_0})| \geq \epsilon - \frac{2M_{\mathcal{F}}}{N}$ and the pair $(f_+(t_{i_0}), f_-(t_{i_0}))$ is the same. This implies that there is a quintuplet $(i_0, s_+, s_-, \bar{\mathcal{F}}_+, \bar{\mathcal{F}}_-)$ such that $i_0 \in \llbracket 1, n \rrbracket$, $(s_+, s_-) \in \mathcal{S}^2$ with $s_+ - s_- \geq \epsilon - \frac{2M_{\mathcal{F}}}{N}$, $\bar{\mathcal{F}}_+$ and $\bar{\mathcal{F}}_-$ are two subsets of $\bar{\mathcal{F}}$ of cardinality at least $\lceil \frac{r}{N^{p+2}} \rceil$, and for each $(f_+, f_-) \in \bar{\mathcal{F}}_+ \times \bar{\mathcal{F}}_-$, the ordered pairs $(f_+(t_{i_0}), f_-(t_{i_0}))$ and (s_+, s_-) are identical. Obviously, any two functions in $\bar{\mathcal{F}}_+$ are ϵ -separated in the metric d_{p, \mathbf{t}_n} , and the same holds true for $\bar{\mathcal{F}}_-$. So, by definition, both $\bar{\mathcal{F}}_+$ and $\bar{\mathcal{F}}_-$ $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N} \right)$ -shatter at least $\text{shat} \left(\lceil \frac{r}{N^{p+2}} \rceil \right)$ pairs. Neither $\bar{\mathcal{F}}_+$ nor $\bar{\mathcal{F}}_-$ $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N} \right)$ -shatters any pair $(s_{\mathcal{T}^q}, \mathbf{b}_q)$ such that $\{t_{i_0}\} \subset s_{\mathcal{T}^q}$. If the same pair $(s_{\mathcal{T}^q}, \mathbf{b}_q)$ is $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N} \right)$ -shattered by both sets, then the pair $(s'_{\mathcal{T}^{q+1}}, \mathbf{b}'_{q+1})$ where $s'_{\mathcal{T}^{q+1}} = \{t_{i_0}\} \cup s_{\mathcal{T}^q}$ and $\mathbf{b}'_{q+1} = (b_0 \ \mathbf{b}_q^T)^T$ is the vector deduced from \mathbf{b}_q by adding one component b_0 corresponding to

the point t_{i_0} , component equal to $\frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (s_+ + s_-) \right\rfloor$, is $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered by $\bar{\mathcal{F}}$. Indeed,

$$\begin{aligned}
\forall f_+ \in \bar{\mathcal{F}}_+, \quad f_+(t_{i_0}) - \frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (s_+ + s_-) \right\rfloor &= s_+ - \frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (s_+ + s_-) \right\rfloor \\
&\geq s_+ - \frac{1}{2} (s_+ + s_-) \\
&\geq \frac{1}{2} (s_+ - s_-) \\
&\geq \frac{1}{2} \left(\epsilon - \frac{2M_{\mathcal{F}}}{N} \right) \\
&\geq \frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}
\end{aligned}$$

and

$$\begin{aligned}
\forall f_- \in \bar{\mathcal{F}}_-, \quad f_-(t_{i_0}) - \frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (s_+ + s_-) \right\rfloor &= s_- - \frac{2M_{\mathcal{F}}}{N} \left\lfloor \frac{N}{4M_{\mathcal{F}}} (s_+ + s_-) \right\rfloor \\
&\leq \frac{1}{2} (s_- - s_+) + \frac{2M_{\mathcal{F}}}{N} \\
&\leq \frac{1}{2} \left(\frac{2M_{\mathcal{F}}}{N} - \epsilon \right) + \frac{2M_{\mathcal{F}}}{N} \\
&\leq - \left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N} \right).
\end{aligned}$$

Summarizing, for each pair $(s_{\mathcal{T}^q}, \mathbf{b}_q)$ $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered by both $\bar{\mathcal{F}}_+$ and $\bar{\mathcal{F}}_-$, we can exhibit by means of an injective mapping a pair $(s'_{\mathcal{T}^{q+1}}, \mathbf{b}'_{q+1})$ $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered by $\bar{\mathcal{F}}$ but not by $\bar{\mathcal{F}}_+$ or $\bar{\mathcal{F}}_-$, so that the number of pairs $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered by $\bar{\mathcal{F}}$ is superior to the sum of the number of pairs $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered by $\bar{\mathcal{F}}_+$ and the number of pairs $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered by $\bar{\mathcal{F}}_-$. This implies that $\text{shat}(r) \geq 2 \cdot \text{shat}(\lceil \frac{r}{N^{p+2}} \rceil)$. Since it has been proved that $\text{shat}(2) \geq 1$, by induction, for all $u \in \mathbb{N}$ satisfying $2N^{(p+2)u} \leq \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n})$, $\text{shat}(2N^{(p+2)u}) \geq 2^u$. Suppose now that we can set $u = \lceil \log_2(K) \rceil$. We then obtain

$$\begin{aligned}
\text{shat}\left(2N^{(p+2)\lceil \log_2(K) \rceil}\right) &\geq 2^{\lceil \log_2(K) \rceil} \\
&> K - 1.
\end{aligned}$$

However, the number of pairs $(s_{\mathcal{T}^q}, \mathbf{b}_q)$ that can be $\left(\frac{\epsilon}{2} - \frac{3M_{\mathcal{F}}}{N}\right)$ -shattered is trivially bounded from above by the total number of those pairs, i.e., $K - 1$. We have thus established by contradiction that

$$2N^{(p+2)\lceil \log_2(K) \rceil} > \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}). \quad (41)$$

A well-known computation (see for instance the proof of Corollary 3.3 in [30]) produces the following upper bound on K :

$$K \leq \left(\frac{e(N-1)n}{d} \right)^d. \quad (42)$$

Substituting (42) into (41) gives:

$$\begin{aligned} \mathcal{M}(\epsilon, \mathcal{F}, d_{p, \mathbf{t}_n}) &< 2N^{(p+2) \left\lceil \log_2 \left[\left(\frac{e(N-1)n}{d} \right)^d \right] \right\rceil} \\ &< 2N^{(p+2) \log_2 \left[2 \left(\frac{e(N-1)n}{d} \right)^d \right]} \\ &< 2 \left[2 \left(\frac{e(N-1)n}{d} \right)^d \right]^{(p+2) \log_2(N)}, \end{aligned}$$

and the last inequality is precisely Inequality (40). ■

References

- [1] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. *Journal of the ACM*, 44(4):615–631, 1997.
- [2] M. Anthony and P.L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, Cambridge, 1999.
- [3] J.-Y. Audibert and O. Bousquet. Combining PAC-Bayesian and generic chaining bounds. *Journal of Machine Learning Research*, 8:863–889, 2007.
- [4] P.L. Bartlett. The sample complexity of pattern classification with neural networks: The size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536, 1998.
- [5] P.L. Bartlett, S.R. Kulkarni, and S.E. Posner. Covering numbers for real-valued function classes. *IEEE Transactions on Information Theory*, 43(5):1721–1724, 1997.
- [6] P.L. Bartlett and P.M. Long. More theorems about scale-sensitive dimensions and learning. In *COLT’95*, pages 392–401, 1995.
- [7] P.L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.

- [8] P.L. Bartlett and J. Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. In B. Schölkopf, C.J.C. Burges, and A. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*, chapter 4, pages 43–54. The MIT Press, Cambridge, MA, 1999.
- [9] S. Bernstein. *The Theory of Probabilities*. Gastehizdat Publishing House, Moscow, 1946.
- [10] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities, a Nonasymptotic Theory of Independence*. Oxford University Press, Oxford, 2013.
- [11] C. Cortes and V.N. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [12] L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag, New York, 1996.
- [13] Ü. Doğan, T. Glasmachers, and C. Igel. A unified view on multi-class support vector classification. *Journal of Machine Learning Research*, 17(45):1–32, 2016.
- [14] R.M. Dudley. The sizes of compact subsets of Hilbert space and continuity of Gaussian processes. *Journal of Functional Analysis*, 1(3):290–330, 1967.
- [15] R.M. Dudley, E. Giné, and J. Zinn. Uniform and universal Glivenko-Cantelli classes. *Journal of Theoretical Probability*, 4(3):485–510, 1991.
- [16] Y. Guermeur. VC theory of large margin multi-category classifiers. *Journal of Machine Learning Research*, 8:2551–2594, 2007.
- [17] Y. Guermeur. A generic model of multi-class support vector machine. *International Journal of Intelligent Information and Database Systems*, 6(6):555–577, 2012.
- [18] M.J. Kearns and R.E. Schapire. Efficient distribution-free learning of probabilistic concepts. *Journal of Computer and System Sciences*, 48(3):464–497, 1994.
- [19] M.J. Kearns, R.E. Schapire, and L.M. Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2-3):115–141, 1994.
- [20] A.N. Kolmogorov and V.M. Tihomirov. ϵ -entropy and ϵ -capacity of sets in functional spaces. *American Mathematical Society Translations, series 2*, 17:277–364, 1961.

- [21] V. Koltchinskii and D. Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1):1–50, 2002.
- [22] A. Kontorovich and R. Weiss. Maximum margin muliclass nearest neighbors. In *ICML’14*, 2014.
- [23] V. Kuznetsov, M. Mohri, and U. Syed. Multi-class deep boosting. In *NIPS 27*, pages 2501–2509, 2014.
- [24] Y. Lei, Ü. Doğan, A. Binder, and M. Kloft. Multi-class SVMs: From tighter data-dependent generalization bounds to novel algorithms. In *NIPS 28*, pages 2026–2034, 2015.
- [25] P. Massart. Some applications of concentration inequalities to statistics. *Annales de la Faculté des Sciences de Toulouse 6^e série*, 9(2):245–303, 2000.
- [26] P. Massart. *Concentration Inequalities and Model Selection: Ecole d’Eté de Probabilités de Saint-Flour XXXIII - 2003*. Springer-Verlag, Berlin, 2007.
- [27] S. Mendelson. Geometric methods in the analysis of Glivenko-Cantelli classes. In *COLT’01*, pages 256–272, 2001.
- [28] S. Mendelson. Rademacher averages and phase transitions in Glivenko-Cantelli classes. *IEEE Transactions on Information Theory*, 48(1):251–263, 2002.
- [29] S. Mendelson and R. Vershynin. Entropy and the combinatorial dimension. *Inventiones mathematicae*, 152:37–55, 2003.
- [30] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. The MIT Press, Cambridge, MA, 2012.
- [31] M.D. Richard and R.P. Lippmann. Neural network classifiers estimate Bayesian *a posteriori* probabilities. *Neural Computation*, 3(4):461–483, 1991.
- [32] M. Talagrand. *The Generic Chaining: Upper and Lower Bounds of Stochastic Processes*. Springer-Verlag, Berlin, 2005.
- [33] V.N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, Inc., New York, 1998.
- [34] V.N. Vapnik and A.Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, XVI(2):264–280, 1971.

- [35] R.C. Williamson, A.J. Smola, and B. Schölkopf. Generalization performance of regularization networks and support vector machines *via* entropy numbers of compact operators. *IEEE Transactions on Information Theory*, 47(6):2516–2532, 2001.
- [36] T. Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251, 2004.